

An Implementation of Process Swapping in MINIX
(A Message Passing Oriented Operating System)

by

Stanley George Kobylanski

B.S., Pennsylvania State University, 1972

A MASTER'S REPORT

submitted in partial fulfillment of the

requirements for the degree

MASTER OF SCIENCE

Department of Computing and Information Sciences

KANSAS STATE UNIVERSITY
Manhattan, Kansas

1989

Approved by:

Mosiah Mizuno
Major Professor

LD
2668
.R4
CMSC
1989
K63
c.2

A11208 317680

CONTENTS

1. INTRODUCTION	1
1.1 Purpose of Operating Systems	1
1.2 MINIX Operating System	2
1.2.1 Introduction	2
1.2.2 MINIX and UNIX	3
1.2.2.1 General Description	3
1.2.2.2 UNIX Description	4
1.2.2.3 MINIX Description	6
1.2.3 MINIX Limitations	9
1.2.3.1 MINIX Problem	10
2. REQUIREMENTS	12
2.1 Extend Memory	12
2.2 Maintain Existing Functions	12
2.3 Maintain MINIX Structure	12
2.4 User Administration	13
2.5 Performance	13
3. DETAILED DESIGN	14
3.1 Introduction	14
3.2 Process Swapping	14
3.2.1 Purpose of Swapping	14
3.2.2 Swapping Functions	15
3.2.3 Alternatives to Swapping	15
3.2.3.1 Increased Memory	15
3.2.3.2 Demand Paging	15
3.3 General Design	16
3.3.1 Introduction	16
3.3.2 Overview	18
3.3.3 Swapping Functions	20
3.3.3.1 Swap Device	20
3.3.3.2 Swapping Out Processes	20
3.3.3.3 Swapping In Processes	26
3.4 Detailed Module Description	27
3.4.1 Introduction	28
3.4.2 Kernel - Management and Decisions	28
3.4.2.1 Swap Task	28

3.4.2.1.1	Swap Task States	30
3.4.2.1.2	Swap Task Transitions	30
3.4.2.1.3	Swap Out Algorithm	33
3.4.2.1.4	Swap In Algorithm	35
3.4.2.2	Clock Task	35
3.4.2.3	System Task	36
3.4.3	Memory Manager - Coordination	37
3.4.3.1	Fork	37
3.4.3.2	Exec	38
3.4.3.3	Paused Process Handling	39
3.4.3.4	Wait Process Handling	39
3.4.3.5	Signal Handling	40
3.4.3.6	Freed Memory Handling	40
3.4.3.7	Swap Out Function	41
3.4.3.8	Swap In Function	41
3.4.4	File System - Swap I/O	42
3.4.5	Swap Device	42
4.	IMPLEMENTATION	44
4.1	Introduction	44
4.2	Kernel Code	44
4.2.1	clock.c	44
4.2.1.1	do_clocktick	44
4.2.2	proc.c	45
4.2.2.1	mini_rec	45
4.2.3	swapper.c	45
4.2.3.1	swap_task	45
4.2.3.2	try_to_swin	46
4.2.3.3	pik_insw	47
4.2.3.4	pik_outsw	47
4.2.3.5	set_swapproc	47
4.2.4	system.c	48
4.2.4.1	do_fork	48
4.2.4.2	do_newmap	48
4.2.4.3	do_exec	48
4.2.4.4	do_xit	49
4.2.4.5	do_lock	49
4.2.4.6	inform	49
4.3	MM Code	49

4.3.1	alloc.c	49
4.3.1.1	alloc_mem	49
4.3.1.2	tot_hole	49
4.3.1.3	compact	50
4.3.2	exec.c	50
4.3.2.1	do_exec	50
4.3.2.2	new_mem	52
4.3.2.3	load_seg	52
4.3.3	forkexit.c	52
4.3.3.1	do_fork	52
4.3.3.2	mm_exit	53
4.3.3.3	do_wait	53
4.3.4	main.c	54
4.3.4.1	main	54
4.3.4.2	reply	54
4.3.5	mswap.c	54
4.3.5.1	do_swout	54
4.3.5.2	swapout	55
4.3.5.3	do_swin	55
4.3.5.4	swap_in	57
4.3.6	signal.c	57
4.3.6.1	check_sig	57
4.3.6.2	unpause	57
4.3.6.3	do_pause	57
4.3.6.4	dump_core	58
4.4	FS Code	58
4.4.1	main.c	58
4.4.1.1	fsinit	59
4.4.2	stadir.c	59
4.4.2.1	do_chdir	59
4.4.2.2	change	59
5.	CONCLUSIONS	60
5.1	Results	60
5.2	Improvements and Further Development	61
	BIBLIOGRAPHY	63
	APPENDIX A - MODIFICATIONS TO KERNEL CODE	A-1
	APPENDIX B - MODIFICATIONS TO MEMORY MANAGER CODE	B-1

APPENDIX C - MODIFICATIONS TO FILE SYSTEM CODE	C-1
----------------------------------------------------------	-----

LIST OF FIGURES

Figure 1-1. Architecture of UNIX Operating System	5
Figure 1-2. Architecture of MINIX Operating System	7
Figure 3-1. MINIX Process Swapping Overview	17
Figure 3-2. Original MINIX User Process States	21
Figure 3-3. MINIX User Process States with Swapping	22
Figure 3-4. Illustration of Compaction	27
Figure 3-5. State Diagram for Swap Task	29

LIST OF TABLES

TABLE 5-1. MINIX NON-SWAPPING vs SWAPPING BENCHMARKS	62
--------------------------------------------------------------	----

ACKNOWLEDGEMENTS

My wife, Janine, and my son, Abraham, have given me more love and support within the past year than any man deserves. It is impossible to return so much, but I will try.

I would like to thank Dr. Virgil Wallentine and Dr. Maarten Van Swaay for serving on my committee and for exposing, to me, a small portion of their wealth of knowledge and wisdom.

I would also like to thank my advisor, Dr. Masaaki Mizuno, for sharing his theoretical knowledge, his technical expertise, and his humor. His guidance has kept me within the bounds of the project while allowing me the latitude to explore until I could reach and defend my own conclusions.

Finally, I would like to thank Charles Clouse, AT&T Document Development Organization, for his help in formatting this document.

CHAPTER 1

1. INTRODUCTION

1.1 Purpose of Operating Systems

A computer system consists of hardware and software that cooperate to do useful work. Computer hardware consists of the physical devices that you can see and touch, such as memory circuits, disk drives, tape drives, terminals, printers, light pens, keyboards, etc. They provide most of the key computer resources: storage for data, computing power for processing data, and devices for input and output of data. Another key resource is the data itself.

Computer software, often referred to as computer programs, consists of a sequence of logical instructions that the computer hardware performs to achieve a desired result. Software and hardware, teamed to form a computer system, have the ability to do various functions, from check book balancing to highly complex programs that can mimic some parts of human behavior.

An operating system is a computer program that:

- manages the resources of a computer, and
- provides easy access to complex hardware.

An operating system allows expensive hardware resources to be shared by many users. For

example, a laser printer that provides high quality printing might be too expensive for a single user, but can be affordable when the cost is shared by many users. Hardware resources can also be hard to use because of their complex interfaces. An operating system hides the hardware complexity by converting user requests so that the hardware can accept them.

Hardware provides "raw computing power" and operating systems make this power conveniently available to users.

1.2 MINIX Operating System

1.2.1 Introduction There are many different operating systems, each built to satisfy certain requirements. MINIX is a general purpose, multiprocessing, multiuser, operating system designed to serve as an aid in teaching operating system concepts. To that end, it was designed to:

- be small, so that it is not overwhelming and can be understood by a student.
- provide an outward appearance (user interface) that mimics a popular operating system called version 7 (v7) UNIX¹ (hence its name *Mini uNIX*). MINIX includes most of the system calls (basic operating system commands), features, and supporting programs provided by UNIX.

1. UNIX is a trademark of AT&T.

- be modular to aid in comprehension and to encourage modification.
- run primarily on the IBM-PC and most compatibles but it has been ported to other machines (e.g., ATARI-ST).
- support multiple users. However, the processing power of the IBM-PC (Intel 8088) is limited such that only 1 user can comfortably be supported. On processors more powerful than the IBM-PC, such as the IBM-AT (Intel 80286) or Intel 80386, more than one user can be supported.

1.2.2 MINIX and UNIX

1.2.2.1 General Description A brief description of MINIX and UNIX is provided here, for more information see references [1] and [2]. MINIX and UNIX consist of four major components: process control, memory management, file system, and input/output.

- Process control handles process creation, communication, scheduling, termination, and miscellaneous process services (suspension, resumption, memory growth, etc.). In MINIX and UNIX, this function is provided by the kernel.
- Memory management allocates and deallocates main memory as needed by processes. In systems that provide swapping or paging, it manages the transfer of process images to and from secondary memory as well as allocation and deallocation of secondary memory. In MINIX, this function is provided by the memory manager (MM). The UNIX kernel provides this function.

- File system manages secondary memory for efficient storage and retrieval of user data. In MINIX, this function is provided by the file system (FS). The UNIX kernel provides this function.
- Input and output procedures provide controlled access to peripheral devices such as terminals, tape drives, disk drives, and network devices for user processes. In MINIX and UNIX, these functions are provided by device drivers (also known as tasks in MINIX) that are linked with the kernel.

1.2.2.2 UNIX Description Unix was designed in the early 1970's to support general computer science research and to provide a custom work environment for its creators. It has evolved from a custom work environment to a popular general purpose operating system.

UNIX is an example of a layered operating system. Its high level architecture is shown in Figure 1-1. The rings in the figure represent levels of interaction and privilege. A program in a ring can only interact with or get services from programs in adjacent rings. For example, the cc program cannot directly interact with (or request services from) the kernel.

Only the operating system (kernel) interacts directly with the hardware, providing common services to programs and insulating them from the hardware idiosyncrasies. Programs interact with the kernel by invoking a well defined set of system calls. An advantage of this design is that the rings can be extended as much as the user prefers.

Within the kernel, there is no structure nor information hiding. It consists of a collection of

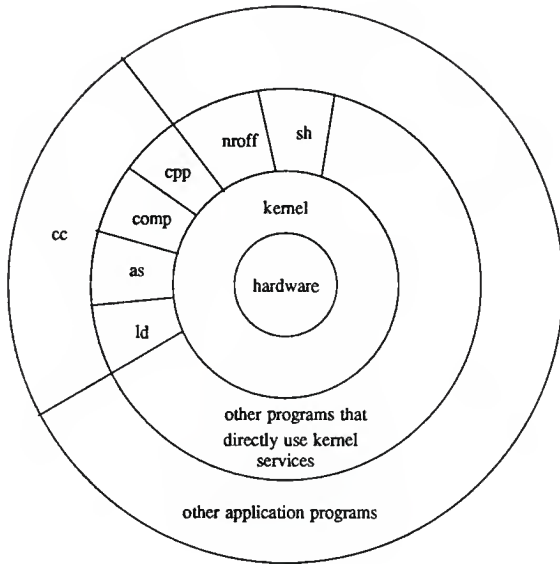


Figure 1-1. Architecture of UNIX Operating System

procedures that are compiled into a single object file. System calls from user processes execute a special instruction called a trap. This instruction switches the machine from user mode to kernel (or supervisor) mode. Control is transferred to the kernel which acts on behalf of the user process. The kernel is not a separate set of processes that run in parallel with user processes, but it is a part of each user process. In effect, the user process becomes the kernel to perform the protected operating system functions in the kernel mode. When the system call completes, the machine is switched back to user mode and the

process resumes in user mode.

As an example, consider the *open* system call which prepares the user process to access a file for reading and/or writing. The user process issues the system call, in the C programming language, as:

```
fd = open(path,mode);
```

where *fd* is the file descriptor used for future file access, *path* is the path and name of the file to be opened, and *mode* is the access permission: read, write, or both. The *open* system call has an entry point in the system call library. The library, encoded in assembly language, contains special trap instructions, which when executed cause an "interrupt" that results in a hardware switch from user mode to kernel mode. For each user process in UNIX, there exists a user stack for use in user mode and a kernel stack for use in kernel mode. The switch to kernel mode causes a switch to the kernel stack and allows the user process to execute the kernel procedures for opening a file. Without getting into great detail about the UNIX file system, the file, if found, is checked for access permissions for this user, prepared for access, and the file descriptor is returned to the system call library *open* routine. The system call library executes another trap instruction that results in a hardware switch back to the user mode. The file descriptor is returned to the user process which resumes in user mode.

1.2.2.3 MINIX Description MINIX is an example of a client-server system. Although it provides the version 7 UNIX user interface, the internal structure and operation are different

than UNIX. In a client-server (or message passing) model, such as MINIX, as much functionality as possible is removed from the operating system leaving a minimal kernel. Most of the operating system functions are contained in user processes. To request a service (system call), a user process (known as a client process), sends the request to a server process which does the work and returns the results. The kernel handles the communication between the clients and servers. The servers run as user mode processes and do not have access to the hardware. The structure of MINIX is shown in Figure 1-2.

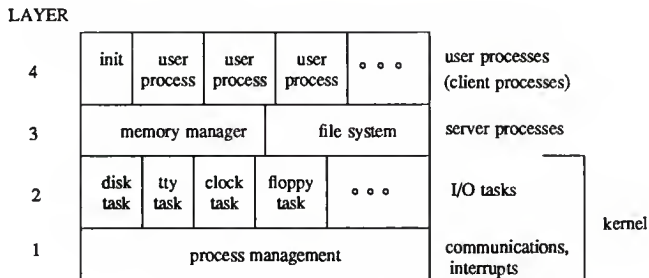


Figure 1-2. Architecture of MINIX Operating System

The layers are similar to the UNIX architecture in that a program can only interact with adjacent layers.

Layers 1 and 2 are compiled into 1 program called the kernel. The process management layer handles all interrupts and traps, and provides message communication between all processes. Layer 2 consists of input/output processes, typically called tasks or device

drivers. Although linked together, they run as separate processes. They provide the device dependent services (physical I/O with specific hardware devices).

Layer 3 contains two processes that support all system calls for user processes. Layer 3 processes provide the device independent services (logical I/O common to all devices).

Layer 4 contains all user processes. System calls from user processes are handled as a client-server transaction. A user process requests a system service by SENDING the request to the appropriate server (FS or MM) and waiting for the results. The server RECEIVES the request, performs the service (which might involve communications with other servers) and returns the results. Then the user process resumes.

For comparison with UNIX operation, a MINIX *open* system call is described. A user process executes the *open* system call using the same syntax and interface as UNIX. The system call library function converts the call to SEND and RECEIVE statements. These are MINIX communication primitives. The library function then executes a trap instruction that causes a software interrupt that is serviced by the kernel. The MINIX kernel puts the SEND and RECEIVE requests on an internal message queue for FS. The user process (via the library function) is then blocked until both requests complete. When FS issues a RECEIVE request to accept new work, the MINIX kernel delivers the user process message to FS. FS copies the message to its data region, determines that it is an *open* system call, and begins to perform the function. All data relating to the file system is maintained within FS. It is not directly accessible by any other process, including the kernel. When FS completes the function, the results are returned using the SEND request with the user

process as the destination. The user process (via the library function) blocked on RECEIVE, accepts the results immediately, freeing the FS SEND block. FS then issues a RECEIVE request for ANYone and awaits new work. When the user process is scheduled to use the CPU, it proceeds from the reception of the system call results.

Note the difference in system call handling between UNIX and MINIX. In UNIX, the user process performs the system call within the boundaries of the kernel procedure. In MINIX, system calls are performed by independent server processes for the waiting user process (client). This has implications on the design presented later in this document.

An empirical study of message oriented operating systems versus procedure oriented operating systems is provided in reference [3]. MINIX fits the definition of a message oriented system and UNIX somewhat fits the definition of a procedure oriented system. The study demonstrates that those two types of operating systems are duals of each other and that a system constructed according to one model has a direct counterpart in the other. It concludes that neither model is inherently preferable, the main consideration being the nature of the machine architecture, not the application that the system will support.

1.2.3 MINIX Limitations MINIX and its host machine(s) do have some limitations. Some of the notable limitations are:

- **Program Size.** MINIX programs are limited to 64K bytes if compiled as non-separate text and data and 128K bytes if compiled as separate text and data. In the latter case, the text and data regions each have a limitation of 64K bytes. The program size

limitation is due to the hardware memory management unit of the host computer.

- **Memory Size.** The main memory size for a MINIX system is limited to 640K bytes. Again, this is a hardware imposed limitation.

1.2.3.1 MINIX Problem A user session with MINIX occasionally results in frustration because of the main memory size limitation. Main memory use on an IBM-PC running MINIX varies, but a typical distribution is:

MEMORY LOCATION	RESIDENT PROGRAMS/DATA
0 - 120K	MINIX Operating System
120K - 360K	root File System RAM Disk
360K - 640K	User Process Area

The amount of main memory available for user processes is about 280K bytes. This allows about four 64K bytes processes to run concurrently. Additional processes will be rejected due to lack of main memory, even though the currently running processes might be blocked and ample CPU power is available. On a multi-process system such as MINIX, it is not uncommon to have many processes active, even for a single user. For example, the user might type:

```
$ ls -al | grep "[a-m]*.c" | more
```

to list specific files in the current directory and view them a screenful at a time. The count of active processes is: user command shell, a shell spun off to handle the multiple

commands line, ls command, grep command, and the more command, for a total of 5 processes. In this example, the user command shell is blocked until all commands complete; it can be swapped out if necessary. The design and implementation discussed in this paper address the memory limitation problem.

CHAPTER 2

2. REQUIREMENTS

This chapter presents the goals and objectives that the design and implementation must meet to provide a useful enhancement to the MINIX operating system.

2.1 Extend Memory

The 1.3 version of MINIX provides about 280K bytes of main memory (about 520K bytes on an AT model) for user programs. Occasionally, this limit is reached without exceeding the processing capability of the CPU. This results in a refusal to perform the user requested command and causes user disappointment and frustration. The system should provide a means to overcome the main memory limitation and allow full use of the CPU. This will allow more useful work to be done and will promote user satisfaction.

2.2 Maintain Existing Functions

Operating system modifications should not affect the current functions provided by MINIX. Additional functions are permissible only if they do not affect the current functions.

2.3 Maintain MINIX Structure

As described in Chapter 1, MINIX is a highly structured, modular operating system. The design should maintain that structure to allow future enhancements to be developed with the same philosophy.

2.4 User Administration

The installation and maintenance of the modification should be minimal and straightforward. Automated installation procedures should be provided. Maintenance procedures should be clearly documented.

2.5 Performance

The implementation should incur minimal performance degradation and main memory usage so as not to offset the performance gain. When the new functionality is not in use, the system performance should equal that of the current system. When in light use, the system performance degradation should not exceed 25%. When in heavy use, the CPU limitations might be exceeded; this cannot be avoided.

CHAPTER 3

3. DETAILED DESIGN

3.1 Introduction

This chapter presents and defends a design that meets the requirements presented in Chapter

2. A process swapper that shuttles process images between main memory and secondary memory has been chosen to satisfy the requirements. The design is applied to the 1.3 version of MINIX.

3.2 Process Swapping

3.2.1 Purpose of Swapping Process memory images² must reside in main memory to run.

If main memory is not available, a process cannot be created. Process swapping allows more processes to run than can fit into the available main memory. It does this, transparently to the user, by shuttling process images between main memory and secondary memory (called the swap device). Clearly, a process on secondary memory requires longer access time than a process in main memory, so low priority, inactive processes are chosen for temporary swap out. When the swapped out processes are ready to run and main memory is available, they are moved back into main memory.

2. A process memory image consists of the data and instructions that must reside in main memory for the process to run. It can include the process text, data, stack, and the process table slot. Within the context of this document, the process memory image consists of the process text, data, and stack regions which can be swapped out and in.

3.2.2 Swapping Functions The functions provided by swapping are:

- provide for process creation when main memory is not available by creating the new process image on the swap device.
- swap in runnable processes and, if necessary, free main memory by swapping out inactive processes.
- choose appropriate processes for swap in/out.
- manage space on the swap device.

3.2.3 Alternatives to Swapping

3.2.3.1 Increased Memory It is not feasible to increase the amount of main memory beyond 640K bytes on an IBM-PC. The IBM-AT model can provide up to 384K bytes of extended memory. Currently, MINIX can use the extended memory only to support the root file system RAM disk. This frees 240K bytes of main memory for user programs and is an excellent alternative to swapping for that machine.

3.2.3.2 Demand Paging Swapping transfers the entire process image between main memory and the swap device. An alternative to swapping is demand paging which transfers individual memory pages instead of entire processes to and from a secondary device. Demand paging permits greater flexibility in mapping the virtual address space of a process into the physical memory of a machine, usually allowing the size of a process to be greater than the amount of available physical memory and allowing more processes to fit simultaneously in main memory. Unfortunately, the IBM-PC does not provide hardware to

support demand paging.

3.3 General Design

3.3.1 Introduction The swapping design affects 3 parts of the MINIX operating system: the kernel, memory manager (MM), and the file system (FS). An overview of MINIX with the swapping design is shown in Figure 3-1. The figure shows most MINIX components but shows the communications only concerned with swapping.

Recall from chapter 1 that tasks are usually interfaces to hardware devices. They are all linked together in the MINIX kernel, but they each behave as a separate process. This gives them independence and access to kernel variables, including the kernel process table. A new task, the swap task, was added to the kernel to manage the overall swapping function. Unlike other tasks, it performs no hardware interfacing. It communicates with other tasks and the memory manager (MM), using normal MINIX communications primitives, to manage the swapping function.

The swap task relies on the memory management functions of MM. MM is responsible for setup and completion of all swap ins and swap outs. MM initiates fork and exec swaps and, upon command from the swap task, performs swap ins and forced swap outs. MM requests physical I/O for swapping via normal user system calls to the file system (FS).

The physical swapping input/output and swap device management is performed by the file system (FS). The swap device is the area on secondary memory used to store process memory images that have been swapped out of main memory.

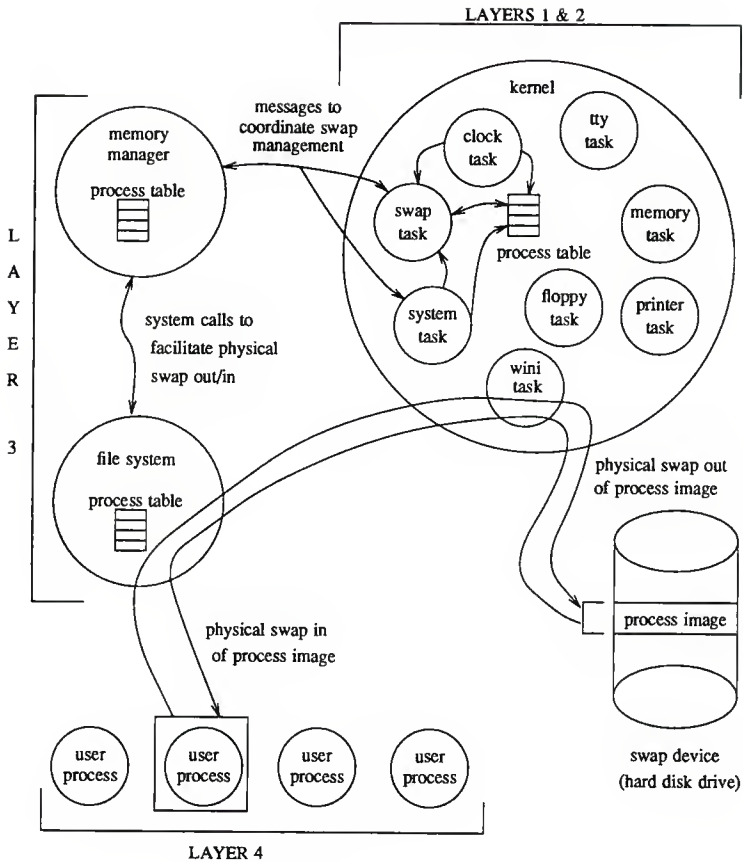


Figure 3-1. MINIX Process Swapping Overview

3.3.2 Overview This part describes a typical situation that involves the entire swapping design (Figure 3-1).

Assume that a user process issues the *fork* system call to create a new process that is a copy of itself (i.e., create a child process). This is typical in multiprocessing type operating systems. The memory manager (MM) receives the request from the user process and starts the *fork* function. *fork* must allocate main memory for the new process image. If main memory is not available, then the new process image is not copied to main memory but is copied to the swap device instead. This is done within MM using normal user system calls to FS. The calling user process memory image is written to a specific directory within a file system on secondary memory (disk). The directory, called the swap device, is a repository for all swapped out process images. When that is complete, MM sends a message (via the system task) to the swap task to indicate that a runnable process is on the swap device. After the swap out MM returns the process id of the child process to the parent process just as in a normal *fork*.

The swap task receives the message and chooses the most eligible, runnable process on the swap device to swap in. It sends a swap in request to MM and awaits the results.

If main memory is available, MM swaps in the process using normal user system calls to FS to read the user process memory image from the directory (swap device) and copy it to main memory. When the swap in has completed, MM notifies the swap task that the process image has been swapped in. If main memory is not available, MM notifies the swap task of the failure. MM also provides information to the swap task to aid in choosing

a process to swap out.

If the swap in succeeded, the swap task checks for more runnable processes on the swap device and attempts to swap them in in the same manner. If the swap in failed, then to free main memory for the runnable process on the swap device, the swap task chooses an eligible process in main memory to swap out. It uses the information provided by MM to aid in selection. If it cannot find a process to swap out, the swap task idles until it receives an indication that main memory has been freed or that a process in main memory might now be eligible to swap out. When a process is chosen for swap out, the swap task sends a swap out request to MM and awaits the results.

MM swaps out the requested process image using normal user system calls to FS to write the user process memory image to the swap device. When the swap out has completed, MM notifies the swap task.

The swap task then chooses the most eligible runnable process on the swap device to swap in. As previously described, it requests MM to swap in the process.

This cycle continues until there are no runnable processes on the swap device. Processes chosen by the swap task for swap out can be runnable or blocked on a system call. Runnable processes are eligible for swap in immediately. Blocked processes are not eligible for swap in until the system call has been completed. When a blocked process on the swap device becomes runnable, a message is sent to the swap task to start the swap in procedure.

3.3.3 Swapping Functions From the Overview, three main swap functions are identified: managing space on the swap device, swapping processes out of main memory, and swapping processes into main memory.

3.3.3.1 Swap Device Swapping introduces another resource that the operating system must manage, the swap device. The swap device is the area on secondary memory used to store process memory images that have been swapped out of main memory. The swap device can be considered an extension of main memory with a large size and limited capability. Its capability is limited to storing process memory images because processes cannot use the CPU while on the swap device. Its advantages are that a process can be kept alive while on the swap device, a system call placed before the process was swapped out can progress (because it is performed by another process: MM or FS), and the main memory that the process is not using can be used by another running process.

3.3.3.2 Swapping Out Processes Figure 3-2 shows the original MINIX user process states. Processes created by *fork* and overlaying programs introduced by *exec* ("new proc") are placed in main memory and then the process is placed on the RUN queue ("on RUN queue"). However, if main memory is not available, then the process creation fails. With the introduction of swapping, these failures become candidates for swap out. Figure 3-3 shows a new state ("SWAPPED & RUNNABLE") to which these failures go when main memory is not available. In this state the process memory images are on the swap device. Later, when main memory is made available, they are swapped in (see "Swapping In Processes").

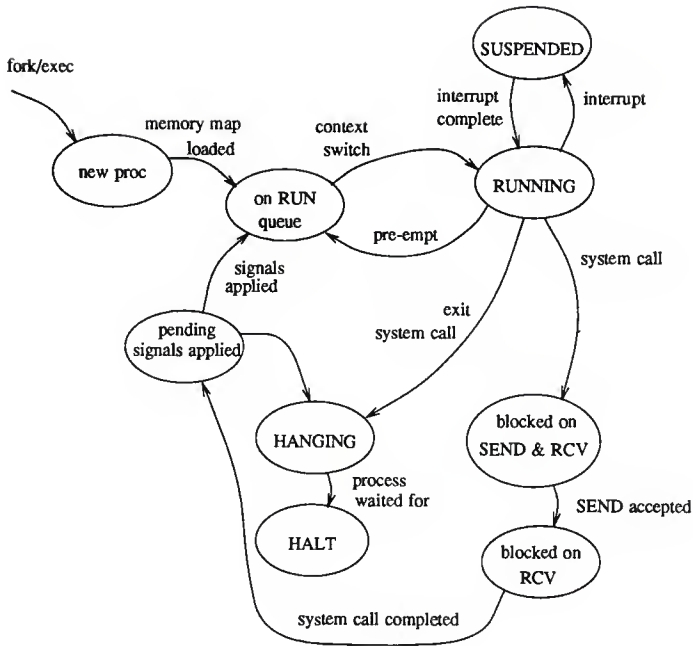


Figure 3-2. Original MINIX User Process States

It is necessary, at times, to force a process out of main memory to free space for a runnable process currently on the swap device. Processes must be chosen carefully for forced swap out. Recall that when a user process performs a system call it is blocked until the system call is completed by MM or FS. Most system calls result in data being transferred to or from the user process data region in main memory. MM and FS rely on the fact that a

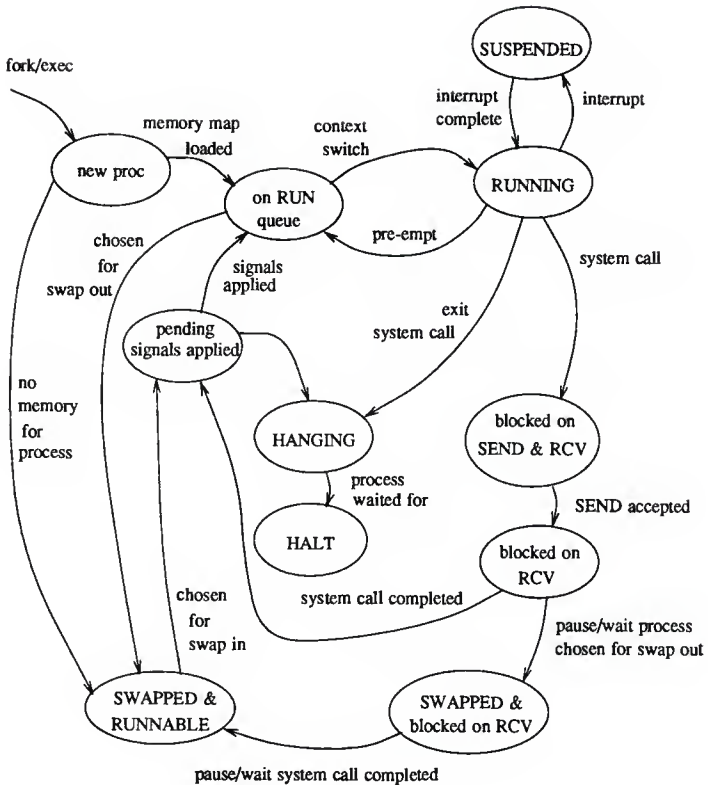


Figure 3-3. MINIX User Process States with Swapping

process will always be in main memory. Therefore, a process that is chosen for swap out should not be waiting for I/O. This affects which processes are eligible to swap out.

For example, suppose that a user process wishes to write some data to an I/O device (state "RUNNING" in Figure 3-2). It issues a *write* system call to the operating system. This is done by SENDING a message to FS and awaiting the result by issuing a RECEIVE from FS. The user process is therefore blocked on SEND and RECEIVE ("blocked on SEND & RCV"). When FS issues a RECEIVE (to accept new work from ANY source), the message passing mechanism of the operating system passes the system call *write* request to FS. The user process is now freed from the SEND block and is blocked only on RECEIVE from FS ("blocked on RCV"). At this point FS is dedicated to processing the user's *write* request. There is no way to stop it. Data will be read from the user's area and copied to the I/O device. If the user process is swapped out BEFORE the data has been retrieved from its data region, FS must wait until the user process is swapped back into core. Currently, FS is not designed to wait for an indeterminate amount of time for that to occur. Similar problems exist with other system calls.

In UNIX, the problem is simplified because the user process itself traps to kernel mode to perform operating system functions. When the user process is in core and in the kernel mode, just before it does an I/O operation, it marks itself as being ineligible for a swap out. It stays marked until the I/O has completed. When the user process is swapped out, no progress can be made on the system call because the user process itself performs the protected kernel procedures in kernel mode. Therefore, an I/O operation cannot begin while the process is swapped out.

Clearly, the UNIX solution cannot be used. In MINIX, one solution is to mark all user

processes that do system calls requiring their data region to be read or written as ineligible for swap out. This requires marking nearly all system calls because the message passing scheme uses memory pointers to the user area to pass system call parameters and to return data and responses. It would not make a significant difference to single out the few that do not access the user process data region.

Another solution is to buffer I/O requests for swapped out processes so that I/O is delayed until the process is swapped in. This is a complex modification to MINIX and the buffering required would increase the operating system size.

This design chooses, for forced swap out, those processes that are blocked for an indefinite amount of time and are expecting a minimal amount of I/O. Conveniently, MINIX clearly identifies such processes in the MM process table. They are processes blocked on the system calls *pause* and *wait*. They are ineligible to run until the system call completes so they are excellent candidates for swap out. This is illustrated in the user process state diagram Figure 3-3 by the new state ("SWAPPED & BLOCKED ON RCV"). These two system calls generate minimal I/O that is easily buffered by the kernel until the process is swapped in. When the system call completes, the ("SWAPPED & RUNNABLE") state is entered and the process is ready for swap in.

A secondary choice for processes to swap out, is processes that are not blocked at all. They are expecting no I/O because they are currently CPU bound. It is an easy matter to swap them out ("SWAPPED & RUNNABLE"). They are not an ideal choice for swap out because they are runnable, but they should rarely get chosen. In situations where they do

get swapped out (no blocked processes remaining in core), they provide a method for sharing the CPU fairly among all runnable processes.

A third possibility for forced swap out is processes blocked on both SEND and RECEIVE. This has been considered and it is feasible, but because of the added complexity, it was not implemented.

Another problem with swapping out processes concerns the location in main memory of the swapped out process and how it relates to the process to be swapped in. MINIX uses the first fit algorithm when allocating main memory for new processes. Suppose that a number of processes reside in main memory with a total free space of 20K bytes. A process on the swap device with a size of 18K bytes is runnable and should be swapped in. However, the free main memory is fragmented and the largest amount of contiguous free main memory is 10K bytes (Figure 3-4a). Should this require a swap out of a program or group of programs whose size is ≥ 18 K bytes, or a swap out of a program(s) next to a free block(s) such that $\text{program(s) size} + \text{free block(s) size} \geq 18$ K bytes? This idea is suggested in references [4] and [5]. One advantage is that no processes are swapped out if a large enough cluster cannot be found. This reduces the number of unnecessary swap outs. (The standard UNIX swapper swaps out processes even though doing so might not result in a large enough area for the incoming process. So unnecessary swapping occurs.)

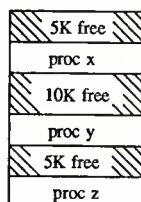
Including this type of requirement in the algorithm for choosing a process for swap out makes it more complex and results in a smaller set of processes from which to choose. Recall that only a special subset of all user processes is eligible for swap out (runnable or

blocked on *pause* and *wait*). An alternative is to compact memory before the swap in such that a contiguous block of free memory of size 20K bytes resides in high memory (Figure 3-4b). The process can then be swapped into the free memory (using first fit) and a potential swap out is eliminated (Figure 3-4c). The performance cost of memory compaction on the IBM-PC is certainly less than the cost of a swap out to secondary memory. In addition, compaction simplifies the algorithm to choose a process to swap out. For example, if an 18K bytes runnable process is on the swap device and total free main memory is 10K bytes and is fragmented, then all that is required to provide space for the swap in is to find an eligible incore program ≥ 8 K bytes, swap it out, and, if necessary, compress memory. The process can then be swapped into the free memory. The overhead of a compaction is incurred but is justified if the following is considered. The amount of main memory available for MINIX user processes is small (about 280K bytes) and the actual number of processes is small (process table can accommodate 15 user processes maximum). Because there is not a lot of main memory to compact, compaction is quick and because there are not a lot of processes from which to choose, the compression algorithm will find an eligible process for swap out much more often than the other algorithm. In addition, the compression algorithm can choose more appropriate processes for swap out.

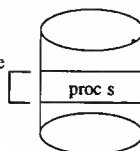
3.3.3.3 Swapping In Processes It now becomes an easy matter to determine which processes are candidates for swap in: processes in the process state ("SWAPPED & RUNNABLE"). A priority, based on residence time, is placed on the processes in that state to facilitate the choice.

MINIX USER PROCESS
MAIN MEMORY AREA

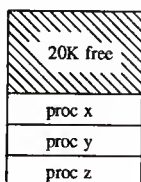
SWAP DEVICE



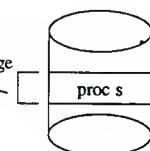
18K bytes
process image



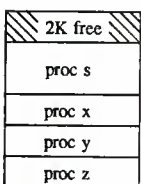
-a- before compaction



18K bytes
process image



-b- after compaction



-c- after swap in

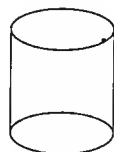


Figure 3-4. Illustration of Compaction

3.4 Detailed Module Description

3.4.1 Introduction The 3 swapping functions just described are provided by various parts of the operating system: kernel provides management and decision making, MM provides coordination, and FS provides primitive swap out/in operations. The design changes to the kernel, MM, and FS to provide swapping are described here.

3.4.2 Kernel - Management and Decisions

3.4.2.1 Swap Task The swap task makes all swapping decisions. The swap task was implemented as a kernel task for the following reasons:

- it can be closely tied to process scheduling
- it needs timely access to some kernel variables.

The swap task, normally idle, awaits messages concerning system swap status and it responds appropriately. Two variables define the state of the swap task: `swap_state` and `inprogress`. `Swap_state` can have the following values:

- `IDLE` - no runnable processes on the swap device
- `SWAP_IN` - runnable process on the swap device.

A runnable process is a process that is not blocked on `SEND` and/or `RECEIVE`. `Inprogress` has the following values:

- `NONE` - no swap operation in progress
- `SWAP_IN` - a swap in is in progress

- SWAP_OUT - a swap out is in progress.

A complete swap task state diagram is shown in Figure 3-5.

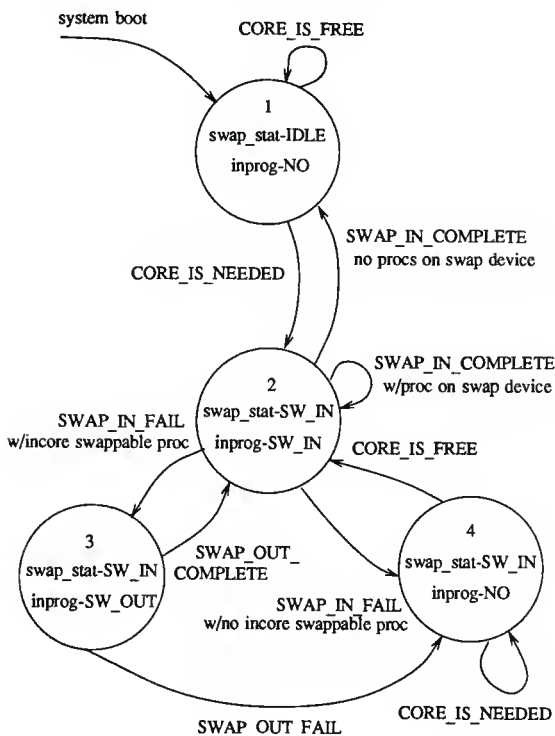


Figure 3-5. State Diagram for Swap Task

The swap task is described in terms of its state and transitions.

3.4.2.1.1 Swap Task States State 1 is entered at system startup. There are no runnable processes on the swap device and the swap task idles whenever in this state.

State 2 is entered when a runnable process is on the swap device and an attempt is being made by the swap task to swap it in.

State 3 is entered when a runnable process is on the swap device, main memory is not available, and a swap out is being attempted by the swap task to free main memory.

State 4 is entered when a runnable process is on the swap device, main memory is not available, and no in core processes are eligible for swap out. The swap task idles until memory is freed or a process becomes eligible for swap out.

3.4.2.1.2 Swap Task Transitions The swap task can receive six different messages from MM and other tasks. The following describes how each of the messages (or events) affects the swap task states and how the swap task responds.

CORE_IS_NEEDED: This message is received when a runnable process is on the swap device. It is the swap task's responsibility to move it into main memory as quickly as possible. This message is sent from:

- MM when a process, on the swap device, that was blocked on the *pause* or *wait* system call has become unblocked.
- MM (via the system task) when a *fork* or *exec* swap has occurred.

If the swap task is in state 1, the swap task's response is to change the swap status from

SWAP_IDLE to SWAP_IN (state 1 to state 2) and try to swap in the runnable process. All runnable processes on the swap device will be evaluated and the most eligible will be chosen (see "Swap In Algorithm" below). The swap task then sends a SWAP_IN_REQUEST message to MM to swap in the chosen process. MM will attempt to swap in the process and will return either a SWAP_IN_COMPLETE or SWAP_IN_FAIL message to the swap task. The swap task can have, at most, one swap in or one swap out operation in progress.

SWAP_IN_COMPLETE: This message is received when a swap in, requested by the swap task (from state 2), has successfully completed. It is sent by MM after swap in completion and should be received only while in state 2. The swap task responds by readying the process for execution. Then the swap task checks for more runnable processes on the swap device (and stays in state 2 or goes to state 1).

SWAP_IN_FAIL: This message is received when a swap in, requested by the swap task (from state 2), has failed. It is sent by MM after attempting a swap in. The reason for failure is assumed to be not enough main memory for the process. MM also sends additional information to help the swap task choose a process for swap out: a list of all processes that are blocked on the *pause* and *wait* system calls and the amount of main memory required to swap in the requested process. Based on the information returned by MM, the swap task responds by choosing a process to swap out so that main memory can be freed to allow the runnable process into main memory. It then marks the chosen process as SWAPPED and sends a SWAP_OUT_REQUEST to MM (state 3). If an eligible process

for swap out cannot be found (state 4), then the swap task sends a `NO_SWAP_OUT` message to the waiting MM and idles until a `CORE_IS_FREE` message is received.

`CORE_IS_FREE`: This message is received when a parameter that controls swap out has changed. If the swap task is in state 4, it can re-evaluate all in core processes to find an eligible swap out candidate. Changeable swapping parameters are: amount of main memory available and process residence time. The other parameter is process memory image size. This message is sent from:

- MM (via the system task) when a process has exited meaning that main memory has become available.
- MM when a process has blocked itself (i.e., it has called the *pause* or *wait* system call), meaning that main memory can be freed by swapping out the blocked process.
- Clock task after the process residence times have increased. This means that a process might now be resident long enough to be eligible for swap out to free main memory.

The swap task (state 4) responds by evaluating all runnable processes on the swap device and trying to swap in the most eligible, hoping that enough core is available (state 2). If the swap in fails, then the swap task will evaluate all processes in main memory for swap out (and go to state 3 or 4).

`SWAP_OUT_COMPLETE`: This message is received when a swap out, requested by the swap task (from state 3), has successfully completed. It is sent by MM after swap out completion. The *fork* and *exec* swaps are reported by the `CORE_IS_NEEDED` message. A

forced swap out is performed only when main memory space is needed for an incoming process. So the reception of this message means that there must be a runnable process on the swap device. The swap task responds by marking the swapped out process as `NO_MAP`, setting the residence time to 0, and marking the process as `BLOCKED` if it is blocked. Then the swap task goes to state 2 and tries to swap in the most eligible process (most likely, the process that caused the forced swap out).

`SWAP_OUT_FAILED`: This message is received when a swap out, requested by the swap task (from state 3), has failed. It is sent by MM to indicate the failure which is assumed to be caused by lack of space on the swap device. The swap task responds by marking the process as not `SWAPPED` and then waiting for the next `CORE_IS_FREE` message (state 4).

3.4.2.1.3 Swap Out Algorithm Processes for swap out are chosen by the swap task based on the following attributes:

- **Size.** The process to be swapped out should be equal to or larger than the amount of main memory needed for the swap in less the current amount of free main memory. This reduces the number of swap outs that would occur if size were not a factor. (Size is not considered in the standard UNIX swap implementations. Studies (see reference [6]) have shown that size considerations can reduce the number of swaps.) When MM returns a `SWAP_IN_FAIL` message, it also sends the size of additional main memory needed for the swap in. The size calculation by MM includes the total amount of existing free memory.

- **Process State.** The process to be swapped out must be blocked on *pause* or *wait* or be runnable. all processes currently blocked on *pause* and *wait* are marked in the MM process table. When MM returns a SWAP_IN_FAIL message, it also sends a list of all processes currently blocked on *pause* and *wait*. These processes have a higher priority for swap out than runnable processes.
- **Time in Core.** Thrashing is a condition in which one or more process images are continually swapped in and out. System resources are expended performing the swapping rather than performing useful work. Thrashing can occur in this design when runnable processes are swapped out because a runnable process on the swap device is immediately available for swap in. If another runnable process is chosen for swap out, then thrashing can occur. To prevent thrashing, a condition for swap out is that a process must reside in main memory for a minimal amount of time. Within that minimal amount of time, the process should get access to the CPU and progress. A potential swap out is delayed due to lack of eligible swap out candidates so a runnable process might spend a bit longer on the swap device. However, system resources are expended doing work more useful than swapping and all users experience better performance.

The following table shows the swap out priority highest to lowest:

SIZE >= NEEDED	BLOCK TYPE	CORE RESIDENCE
Y	pause/wait	oldest in this category
Y	none	oldest above minimal
N	pause/wait	oldest in this category
N	none	oldest above minimal

3.4.2.1.4 Swap In Algorithm Only runnable processes are chosen for swap in. The runnable process that has been on the swap device the longest is chosen for swap in. No minimum amount of time on the swap device is required because thrashing is minimized by the swap out algorithm.

3.4.2.2 Clock Task The clock task has 2 swapping related duties:

- Every second, the clock task increments the residence time of each user process. The residence time is the amount of time that the process has recently been in its current residence (main memory or on the swap device). It is set to zero whenever a process is:
 - created,
 - moved into main memory,
 - and moved onto the swap device.

Residence time is used as a parameter in determining the eligibility of a process for swap out. Thrashing is minimized by requiring that a process remain in main memory

for a minimum amount of time before being swapped. Residence time is also an attempt, although weak, to allow equal access to the CPU for all processes. This is especially true for runnable processes that are chosen for swap out.

- The clock task checks the status of the swap task every second. If a runnable process is on the swap device and the swap task is not in a swap in state, the clock task will send a `CORE_IS_FREE` message to the swap task. This means that residence times have changed and that it might now be possible to find an eligible process for swap out. The swap task can then re-evaluate all processes.

3.4.2.3 System Task MM and FS do not have direct access to the kernel variables and kernel process table. The system task is their interface to the kernel to get data, set kernel variables, and perform miscellaneous duties. Some of the existing communications between the system task and FS and MM are used as vehicles for communicating swapping related data. The swapping related duties performed by the system task are:

- When MM performs a *fork*, it notifies the kernel to set up the kernel process table. The swapping design takes advantage of the existing *fork* communication when a *fork* swap occurs. If the *fork* message from MM indicates that a *fork* swap has occurred, the system task marks the kernel process table slot as `SWAPPED`, and, if the swap task is `SWAP_IDLE`, sends a `CORE_IS_NEEDED` message to the swap task.
- When MM performs an *exec*, it notifies the kernel to set up the kernel process table. If the *exec* message from MM indicates that an *exec* swap has occurred, the system task marks the kernel process table slot as `SWAPPED`. If the swap task is `SWAP_IDLE`, the

system task sends a `CORE_IS_NEEDED` message to the swap task to indicate that a runnable process exists on the swap device. If the *exec* message indicates that a normal *exec* has occurred and the swap task status is `SWAP_IN`, a `CORE_IS_FREE` message is sent to the swap task. This means that an *exec* has occurred that might have freed some main memory.

- When a process exits, MM notifies the kernel via the system task to clean up the process's process table slot and free it for reuse. If the swap task status is `SWAP_IN`, a `CORE_IS_FREE` message is sent to the swap task to indicate that main memory has been freed.

3.4.3 Memory Manager - Coordination

3.4.3.1 Fork The *fork* system call creates a new process by copying the current process's text, data, stack, and process table slot. The parent and child receive a different return value to enable them to identify themselves. When main memory is allocated for the child process, the allocation routine determines whether enough memory is available. If enough contiguous free memory is available, then the *fork* can continue. If free memory is available, but is not contiguous, the allocation routine compacts main memory (i.e., relocates the process images in main memory to create a large contiguous free memory space in high memory). Compaction was implemented to serve two purposes:

- reduce the number of swap outs
- simplify the swap in/out algorithms.

In either case, the memory required for the child process is allocated. However, if the required main memory is not available, the *fork* swap out is begun. It calls the *swapout* procedure (see "Swap Out Function") to copy the parent's process image to the swap device. It marks the MM process table as FKSWAPPED (*fork swapped*), notifies the kernel that a *fork* swap has occurred (see "System Task"), and notifies the kernel of the child's memory map for swap in.

3.4.3.2 Exec The MM *exec* system call overlays the existing program with a new program. (The program changes while the process remains.) The new program inherits the environment of the calling program. When main memory is allocated for the new program, the allocation routine determines whether enough memory is available. The memory space of the current process is included in the free memory evaluation. If enough memory is available, then the *exec* can continue. If enough memory is available, but is not contiguous, then the memory of the existing process is freed, main memory is compacted, and the memory required for the new program is allocated.

If the required main memory is not available, the *exec* swap out is begun. The new program text and data are read from the executable file on disk and written to the swap device. The stack is adjusted and also written to the swap file. The existing memory is freed, the MM process slot is marked as SWAPPED, and the kernel is notified that an *exec* swap has occurred (see "System Task").

3.4.3.3 Paused Process Handling There are two cases to consider:

- when any process issues the *pause* system call
- and when a *pause* swapped process becomes unpaused.

When a process issues the *pause* system call, it is blocked until the system call completes and it expects no I/O. This makes it an excellent candidate to swap out when space is needed in main memory. If, at the same time, a runnable process is on the swap device waiting to be swapped in, a `CORE_IS_FREE` message is sent to the swap task. This means that a process is available for swap out.

A process leaves the *pause* state upon reception of a signal. If the process is swapped out, and there are no other runnable processes on the swap device, a `CORE_IS_NEEDED` message is sent to the swap task. The swap task marks the kernel process table as not `BLOCKED`. After the process is swapped in, the expected `EINTR` result value is returned to the calling process.

3.4.3.4 Wait Process Handling There are three cases to consider:

- when a process issues the *wait* system call
- when a waiting, swapped process is unblocked by a signal
- when a waiting, swapped process receives a death of child signal.

The first two cases are the same as those described for *pause* process handling.

The death of child requires a little more handling. Death of child requires that the waiting process receive the exit status of the dead child process. This is accommodated for swapped processes by postponing the process slot cleanup of the dead process, marking it as **HANGING** (i.e., a zombie), freeing the dead child's memory, and saving the dead child's process slot number in the parent's process table slot. When the process is swapped in, the dead child's exit status is returned to the parent and the process slot cleanup is performed.

3.4.3.5 Signal Handling If a signal is sent to a swapped out process it is saved until the process has been swapped in. If the process is **BLOCKED**, it is marked not **BLOCKED** in the kernel process table and a **CORE_IS_NEEDED** message is sent to the swap task which will eventually start the swap in for this process. All signals received are stored in the kernel process table. When the process is swapped in and before it is marked ready to run, all pending signals are applied. Applicable signal actions are:

- terminate the process - the user process is terminated,
- ignore the signal - the signal is ignored,
- and perform a function - the user process stack is adjusted to set up a call to the signal handling function. When the process runs the signal handling function will be the first thing performed.

3.4.3.6 Freed Memory Handling Main core is freed by the *exit* system call and sometimes by the *exec* system call (i.e., execing a program with a smaller process memory image). The system task *sys_xit* and *sys_exec* functions determine if a runnable process is

on the swap device when either of the above occurs. If so, then they send a `CORE_IS_FREE` message to the swap task.

3.4.3.7 Swap Out Function Preparation for physical swapping of a process image to the swap device is performed by MM. This occurs two ways, *fork* swap and forced swap out. (*exec* swaps are a little different because the process image is in the executable file, not in main memory.) The swap task sends a `SWAP_OUT_REQUEST` message to the `do_swap_out()` function. The process is marked as `SWAPPED`. The location and size of the process segments are determined. A file name for the process is generated. The file is created and the process segments are written to the swap device. If the request is from the swap task, main memory is freed and a `SWAP_OUT_COMPLETE` message is sent to the swap task.

3.4.3.8 Swap In Function Preparation for physical swapping of a process image from the swap device is performed by MM. The swap task sends a `SWAP_IN_REQUEST` message to the `do_swap_in()` function which determines if there is enough main memory for the incoming process. If main memory is not available, then the function makes a list of all processes that are blocked on the *pause* and *wait* system calls. It sends this list along with the size of main memory needed and the `SWAP_IN_FAILED` message to the swap task. This aids the swap task in choosing a process for swap out.

If there is enough free main memory, then the swap in proceeds. If enough free memory exists but is not contiguous, then the allocation routine compacts memory, as described for

fork and *exec*, and the swap in proceeds. The memory location and size of each process segment are determined. The file name for the process is re-generated. The file is opened and the process segments are read from the swap device and copied to main memory. The process is marked as not **SWAPPED** and the swap file is removed. If the process was swapped out as a result of a *fork*, then a wakeup message is sent to the process just as in a normal *fork*. All signals that were received by the process while swapped out are applied. If the process was *pause* or *wait* and was interrupted by a signal, then the **EINTR** result is sent to the process. If the process was *wait* and was awakened by the death of a child, then the child's exit status is sent to the process and the dead process's process table slot is cleaned up. The **SWAP_IN_COMPLETE** message is then sent to the swap task.

3.4.4 File System - Swap I/O

3.4.5 Swap Device The swap device is the area on secondary memory that is used to store process images that have been swapped out of main memory. Swap device input/output operations should be efficient to minimize system performance degradation. Typically, space on a swap device is allocated in multiblock segments so that the swapped process images are written and read contiguously. Contiguous, multiblock reads and writes are quicker than normal file system block-by-block operations. When a process image is swapped into main memory, its storage space on the swap device is freed. Space management routines should be used to manage fragmentation that occurs on swap devices. Because the allocation scheme for the swap device differs from the normal file allocation scheme, data structures that catalog free space differ too.

The above requirements should ensure good swap device input/output performance but they also exact a price: additional code and complexity to perform and maintain that function. On a system, such as MINIX, where the main memory for user processes is limited, additional code in the operating system will further limit that memory.

This design uses an approach which sacrifices efficiency for simplicity. The existing file system is used to manage the swapping input/output operations. The swap device, therefore, consists of a normal file system directory that is protected (with normal file system protection mechanisms) from all users except root. The file system handles all swapping input/output as it handles all other file system input/output. No special contiguous, multiblock input/output operations or swap device space management procedures are required. Normal *creat*, *open*, *read*, *write*, *seek*, *close*, and *unlink* system calls are used. The advantages are: simple implementation, simple administration, and no large operations overhead. Disadvantages are: slow swapping input/output operations and conflicts with the file system. (When the file system is swapping, it is not free to do useful user program input/output.)

The file system manages the swap device as it does other regular files. The only change to the file system allows MM to go to the swap directory directly. The inode of the swap directory is determined and saved by the first access to the swap directory. Subsequent accesses use the saved information.

CHAPTER 4

4. IMPLEMENTATION

4.1 Introduction

This description of the MINIX swapping implementation shows how the design elements presented in Chapter 3 are implemented. It assumes a knowledge of the MINIX version 1.3 operating system and the design described in Chapter 3. All swapping code is written in the C programming language just as most of MINIX. The code itself is commented, so the descriptions here will be brief. The description is organized into 3 parts: kernel, MM, and FS. The code is listed in the appendix³. New programs and functions are listed in their entirety. All C functions that have been altered are listed in their entirety. Only changes to header (.h) files are listed.

4.2 Kernel Code

4.2.1 clock.c

4.2.1.1 do_clocktick Lines 64 through 69 update the kernel process table residence times for each active process by a simple increment. Residence times are reset to 0 when a process is created, moved to the swap device, and swapped into main memory. The

3. The original portions of the code, copyrighted by A. S. Tanenbaum, are reprinted with his permission with the restriction that reproduction be limited.

residence time variable is an unsigned integer with a maximum value of 65535. The value will turnover to 0 after about 18 hours; this is not considered to be a problem.

When there is a runnable process on the swap device, lines 70 through 75 send a `CORE_IS_FREE` message to the swap task as a notification that process residence times have changed. The swap task might now be able to choose a process for swap out.

4.2.2 `proc.c`

4.2.2.1 `mini_rec` Lines 53 and 54 were changed to check for a pending message from the swap task to MM. See discussion about non-blocking message transfer in "Kernel Code, `swapper.c`". If MM is waiting for a message, and a message from the swap task is ready, then the `inform()` function is called to deliver the message.

4.2.3 `swapper.c`

4.2.3.1 `swap_task` Upon entry, the `init_swap()` function is called to set the initial state for the swap task. Then, as with other tasks, a large case statement is entered to handle all messages that can be received.

`CORE_IS_FREE` (line 88) corresponds to the detailed design description. The `try_to_swin()` function checks for runnable processes on the swap device, picks the most eligible, and sends a swap in request to MM.

`CORE_IS_NEEDED` (line 94) first determines whether the message indicates that a swapped process has become unblocked. If so, the kernel process table slot is marked as

not BLOCKED. If swap_stat is SWAP_IDLE, then a swap in is attempted.

SWAP_IN_COMPL (line 105) indicates that the requested swap in has successfully completed. The process table is updated and the process is made ready to run. The swap device is checked for more runnable processes.

SWAP_IN_FAILED (line 118) indicates that the requested swap in has failed. MM has sent the size of core needed and a list of all processes blocked on the pause() and wait() system calls. MM is blocked awaiting a reply to this message. the pik_outsw() function is called to choose a process for swap out and a message is sent to MM to request the swap out. If no eligible process is chosen for swap out, the swap task responds to MM with NO_SWAP_OUT.

SWAP_OUT_COMPL (line 139) indicates that the requested swap out has completed. The NO_MAP flag is set and the BLOCKED flag is maintained by the kernel to indicate whether a swapped process is blocked or runnable (not blocked).

SWAP_OUT_FAILED (line 152) indicates that the requested swap out has failed. No action is performed until another message is received.

4.2.3.2 try_to_swin This function checks for runnable processes on the swap device (process table flags of SWAPPED and BLOCKED), chooses the most eligible, and requests MM to perform the swap in.

4.2.3.3 pik_insw This function returns the process slot number of the oldest, runnable process on the swap device.

4.2.3.4 pik_outsw This function returns the process slot number of the most eligible process to swap out. It does this by assigning a value to each process based on a number of process attributes. An initial priority value of 0 is assigned to each process. If the process is at least as large as the size needed, then PICKSIZE (10000) is added to the priority value. If it is blocked on the pause() or wait() system call and has been resident for the minimal interval CRESMIN (15 seconds), then PICKPWBLK (100) is assigned. If it is not blocked at all and has been resident for CRESMIN, then PICKNBLK (1) is added. If it is ineligible for swap out, PICKINELL (-20000) is added. The process with the highest priority value above 0 is returned.

4.2.3.5 set_swapproc A task is a high priority process and MM is a medium priority process. A process should send a message directly to a lower priority process only when it is known that the lower priority process is waiting for a message. It is not always known what MM is doing, so a non-blocking message transfer has been devised.

To send a non-blocking message to MM, such as for a SWAP_IN_REQUEST, a special kernel global structure was created called swap_proc. It contains a status which can be either MPENDING (a message is waiting to be sent to MM) or NOMES, and a message pointer. This function builds the message, sets the status to MPENDING, and sets the message pointer. The inform() function is called. If MM is waiting, the message will be

delivered immediately. If not, then it is delayed, but the swap task is not blocked. Whenever MM requests a RECEIVE for ANY, the mini_rec() function (in proc.c) checks for a non-blocking message from the kernel and the message is eventually delivered.

4.2.4 system.c

4.2.4.1 do_fork Line 39 sets residence times for all forked processes to 0, whether swapped or not.

Lines 42 through 50 are performed when a fork swap has occurred. The process table is marked and if the swap task is idle, a CORE_IS_NEEDED message is sent. The PROC1 message variable is set to 0 to indicate that the process is not blocked.

4.2.4.2 do_newmap When an exec swap occurs, the main memory space of the execing process is used as a data buffer. Data is copied from the executable file to this main memory area and then written to the swap file. (See MM exec.c for more details.) To set up this memory area as a buffer, exec changes the memory map of the execing process by calling the sys_newmap() function. Lines 99 through 105 check for this situation and set the process table flag to NOMAP to indicate that.

4.2.4.3 do_exec If an exec swap has occurred, the exswap flag is set. Lines 133 through 140 set the process table to SWAPPED and send a CORE_IS_NEEDED message to the swap task if it is SWAP_IDLE.

If an exec has occurred without a swap out and the swap task status is SWAP_IN, lines 144

through 146 send a `CORE_IS_FREE` message to the swap task.

4.2.4.4 do_xit When a process dies, main memory is freed. If the swap task status is `SWAP_IN`, then lines 212 through 215 send a `CORE_IS_FREE` message to the swap task.

4.2.4.5 do_lock This is a new function that allows MM to turn on/off hardware interrupts. It is used when MM is performing memory compaction.

4.2.4.6 inform If a non-blocking message from the swap task to MM is queued, then lines 254 through 259 will send the message immediately (see `set_swapproc`).

4.3 MM Code

4.3.1 alloc.c

4.3.1.1 alloc_mem This function has been changed, such that, if a request for contiguous memory cannot be satisfied, but enough fragmented free memory exists, then the entire memory area is compacted to satisfy the request.

4.3.1.2 tot_hole This is a new function that calculates the sum of all free main memory. It replaces the `max_hole()` function which returned the largest hole of free memory. Because compaction is used, the total free memory size is more useful than just the largest hole size.

4.3.1.3 compact This is a new function that performs the memory compaction. It begins by calling the `sys_lock()` function to turn off all system hardware interrupts so that no I/O transfers are done while a process image is being moved to another memory location. A loop is entered where the first free hole is found. The size of the process just above the hole is calculated and then the process image is copied into the hole. The process' original memory is freed and core is allocated for the process' new location. The process table memory map is adjusted. The loop repeats until there is only 1 hole remaining. The hardware interrupts are then enabled.

4.3.2 exec.c

4.3.2.1 do_exec This function has been changed to add the exec swap out feature. The call to the `new_mem()` function on line 63 returns `EXSWAPD` if main memory is not available for the execed program. `New_mem()` also saves a copy of the existing memory map for this process.

If a swap out is necessary, lines 93 through 185 are performed. First the swap file is created by changing to the swap directory, checking access permissions, and creating the swap file. Any failures up to this point result in an `ERROR` return value which indicates that the program could not be execed. The user process can handle that error any way that it chooses. Failures beyond this point result in system errors because the execing process is destroyed.

To put the execed program onto the swap device, the entire execing process memory area

(up to 2047 blocks) is converted to a temporary data buffer. The `sys_newmap()` function is called to notify the kernel of the change. The third parameter, `TRUE`, indicates that this will be a data buffer and that the system task should mark the kernel process table with the `NO_MAP` flag to prevent this process from being chosen to run.

The executable program is then copied to the swap device within the loop beginning at line 136. The text segment is copied into the data buffer of the execing process via the `load_seg()` function and then written to the swap file via the special version of the `write()` system call in which FS copies data directly from the user process memory area to the disk file without going through MM. After the segment has been completely copied, the file size is increased so that it is equal to the length of the segment as indicated in the MM process table memory map. This makes swap in much easier. The loop is repeated for the data segment.

The stack is created from the original stack and written to the swap file in the same manner as above.

The memory area is freed and the MM process table is marked as `SWAPPED`. The call to the `sys_exec()` function tells the kernel (swap task) about the exec and possible swap out (see Kernel Code - `system.c` - `do_exec`).

The fork swap out and forced swap out both call the `swapon()` function in `mswap.c`. They copy process images from main memory to a swap file. The exec swap does not call `swapon()` because it must copy the program from the executable file, transform the data

into a process image, and then write it to the swap file. In addition, the amount of memory used by the execing process is used as a buffer for the file transfer and transformation and might not be large enough for the entire new program. So it must be done piecemeal. Rather than add this complexity to the `swapout()` function, it was included here.

4.3.2.2 new_mem This function has been changed, lines 268 through 281, to check for total free memory rather than just largest contiguous free memory hole. It also includes the size of the execing process image in the free memory calculation. If memory is not available, `new_mem()` generates the memory map for the new program in a temporary variable and returns `EXSWAPD` to indicate that an exec swap out should be done.

The size of the execing process image is also calculated. If a swap out is necessary, the size will be used to determine how large the data buffer can be.

4.3.2.3 load_seg The new input parameter "usr" (line 340) allows the calling function to specify the user process whose memory area is being loaded. The previous version defaulted to the execing program. The new parameter is useful for the `swap_in()` function in `mswap.c`.

4.3.3 forkexit.c

4.3.3.1 do_fork This function has been changed to add the fork swap out feature and to use the memory compaction feature to reduce the number of fork swap outs. The check for free memory, lines 38 through 45, looks at total free memory rather than the largest

contiguous free memory area. If memory is not available, then a fork swap out must be done. The memory map is adjusted, if necessary. The user program controls the program stack pointer and can increase the stack size beyond the stack region. MINIX is not aware of it until the program requests more data space. When that occurs, MINIX checks the stack pointer and adjusts the stack region size in the memory map to correspond to the new stack pointer. The `adjust()` function verifies and adjusts, if necessary, that the stack pointer is within the stack region as defined in the memory map. The memory map is used as a measure of the process image size for swap out and swap in.

Lines 77 through 81 call the `swapout()` function to copy the parent's image to the swap file.

Lines 86 and 87 set the new MM process table variable to 0.

Line 109 notifies the kernel (via the system task) as to the status of the fork. If a fork swap has occurred, the system task is not notified of the child's new map nor is the reply returned to the child (because it is swapped out). It will be returned after the child is swapped in.

4.3.3.2 mm_exit This function has been changed, lines 164 through 173, to check for waiting parents that are swapped out. If so, then the cleanup of the process table slot and return of the exit status to the parent is delayed until the parent is swapped in. The memory occupied by the exiting process is freed in either case.

4.3.3.3 do_wait This function has been changed, lines 225 through 234, to check for processes on the swap device. If so, then a `CORE_IS_FREE` message is sent to the swap

task to indicate that the process that has just called wait() is eligible for swap out.

4.3.4 main.c

4.3.4.1 main In the call to the service functions on line 27, a dummy parameter (0) was added. This allows the service functions, specifically do_swout(), to be called with parameters by other functions (see "MM Code, do_swin").

4.3.4.2 reply Line 52 was added to exclude the swap task from the validation check because MM does send replies directly to the swap task and the swap task is not in the MM process table.

4.3.5 mswap.c This new file contains procedures for swap out and swap in.

4.3.5.1 do_swout This function is called from main() when the SWAP_OUT_REQUEST is received from the swap task. It is also called from do_swin() after a swap in fails and the swap task suggests a process for swap out. It starts by getting the number of and pointer to the slot of the process to be swapped out. The swapout() function is called to perform the swap out. If the swap out succeeds, the process slot is marked as SWAPPED, and if the process is blocked on the pause() or wait() system call, the PROC2 message variable is set to TRUE. The SWAP_OUT_COMPLETE message is sent to the swap task. If the swap out failed, the SWAP_OUT_FAILED message is returned.

4.3.5.2 swapout This function is called from the `do_swout()` function for a forced swap out and from the `do_fork()` function for a fork swap out. The big difference between the two is that the fork swap does not want to remove the parent process image from main memory, but the forced swap out does. The `clear_mem` variable is `FALSE` for the former and `TRUE` for the latter. The `dump_core()` function is called to perform the actual swapout. If it fails, an error message is returned to the caller. If it succeeds, and if it is a forced swap out, then the memory of the process is freed.

4.3.5.3 do_swin This function is called from `main()` when the swap task sends a `SWAP_IN_REQUEST` message. `do_swin()` gets the process table slot number of the process to swap in, calculates the amount of memory needed, and tries to allocate it. If memory is available, then `do_swin()` marks the process as not `SWAPPED` and calls the `swap_in()` function to perform the swap in. After swap in completion, many items are checked and some cleanup chores are performed.

If the process was fork swapped, then a reply of 0 is sent to indicate that it is a child of the fork.

If any signals were sent to the process while it was swapped out, then they must be applied to the process. (The actual signal handling will be done when the process is put into the `RUNNING` state.) During swap out, all received signals are stored in the `mp_ssw_map` process table variable. For each signal that was received, the `sig_proc()` function is called to process it.

If a process was blocked on the `wait()` system call and one of its child processes died, then the zombie child can now be laid to rest, its process table slot freed, and its exit status returned to its parent. All done by the `cleanup()` function. The slot number of the HANGING child is saved in the `mp_deadchild` process table variable.

If a process was blocked by a `pause()` or `wait()` system call and a signal(s) was sent to the process to wake it up, then the `EINTR` (i.e., system call interrupted) error value must be returned to the process.

The `SWAP_IN_COMPLETE` message is returned to the swap task.

If memory is not available for the swap in, then this function gathers knowledge known only by MM and sends it to the swap task. With this knowledge, the swap task can make a more informed choice of which process to select for swap out. It gathers the slot numbers of all processes in core that are blocked on the `pause()` and `wait()` system calls and stuffs them bitwise into a 32-bits variable (`LONGI`). It also calculates the amount of main memory needed for the swap in (`PROC1`).

It is necessary to freeze the user processes at this point so that they cannot progress and change their pause/wait status. So the `SWAP_IN_FAIL` message is sent to the swap task using the `send_rec()` function. MM blocks on `RECEIVE` from the swap task. If the swap task cannot choose a process for swap out, then no action is taken. If the swap task chooses a process for swap out, the `do_swout()` function is called to perform the task.

4.3.5.4 swap_in This function is called by the `do_swin()` function. It generates the swap file name from the process slot number and opens the file. The memory map is changed to reflect the new base address in main memory. The gap location is initialized to zeros for security. Each segment, text, data, and stack is copied from the swap file via the `loadseg()` function. At completion, the swap file is removed.

4.3.6 signal.c

4.3.6.1 check_sig Lines 56 through 61 were added to handle swapped out processes to which a signal was sent. The `sig_proc()` function is delayed until the process is swapped in. The fact that a signal was sent is recorded bitwise in the process table variable `mp_ssw_map`. Note that the `unpause()` function is called just after this.

4.3.6.2 unpause Lines 117 through 121 check for swapped processes that are also paused. If so, the process table is marked as `WASPWS` (was pause/wait and a signal was received).

The process is now runnable, so a `CORE_IS_NEEDED` message is sent to the swap task with the process table slot number of the awakened process. The swap task will mark the process table as not `BLOCKED` and attempt to swap in the process. Lines 129 through 133 accomplish the same thing for waiting processes.

4.3.6.3 do_pause Lines 85 through 92 send a `CORE_IS_FREE` message to the swap task to indicate that a process has blocked itself and is available for swap out if necessary.

4.3.6.4 dump_core This function, as its name implies, was used only for writing the process image to a file on disk. It has been altered to serve another purpose: swap a process image out to the swap device. It just so happens that these two functions are similar enough to handle them in a general manner. Changes were made throughout the function so it will be described in its entirety.

The input parameter "type" is 0 for a dump core request, 1 for a swap out, and 2 for a fork swap out. The function changes to the working directory of the user process for dump core or the swap device for swap outs. The file is checked for write access permission. For swap out, the swap file name is generated from the process table slot number. The file is created; if it fails, an error is returned.

The memory map is checked for accuracy (i.e., is the stack pointer within the stack region?) and adjusted, if necessary.

For core dumps, the memory map is written to the core file. Each process segment, text, data, and stack is written to the core file or swap file (as appropriate). The write() system call to FS is special in that the data is copied directly from the user memory area and written to disk without going through the MM.

4.4 FS Code

4.4.1 main.c

4.4.1.1 fsinit A new variable, `swap_node` is declared and initialized to `NIL` on lines 10 and 29. It is used to store the inode of the swap directory.

4.4.2 stadir.c

4.4.2.1 do_chdir Lines 38 through 49 were added to allow MM to quickly change to the swap directory when the `cd_flag` is 2. If the directory has not been opened yet, then the `change()` function is called to get the inode and store it in the `swap_node` variable.

4.4.2.2 change Lines 80 through 83 explicitly set the `user_path` to the `name_ptr` if the `cd_flag` is 2.

CHAPTER 5

5. CONCLUSIONS

5.1 Results

The implementation of swapping satisfies the requirements of Chapter 2:

- It successfully extends the limited memory of the system. The CPU is more fully utilized and user frustration is reduced.
- All existing functions are maintained; no new functionality was added.
- The MINIX structure is maintained; related extensions are suggested below.
- Installation is simple. It requires installing a protected directory "/usr/swap". Maintenance requires ensuring that enough secondary memory is available for the swap device.
- Performance is acceptable. It is not degraded when swapping is not in use and shows acceptable degradation when in light use. Table 5-1 shows benchmarks of system performance under various conditions. The benchmarks are based on compiles of the MINIX kernel, MM, and FS using the MINIX supplied compiler and makefiles. Note that there are no perceptible compile time differences between MINIX version 1.3 without swapping and version 1.3 with swapping when no swapping occurs. When swapping occurs, the performance starts to suffer. The AT class machine takes about 23% more time to concurrently compile the kernel and MM while the XT class only

requires about 11% more time. This is considered light to medium swapping and because the compiles are CPU intensive, it is acceptable. For heavier swapping, the AT class machine takes about 52% more time to concurrently compile the kernel, MM, and FS. This is probably exceeding the capacity of the CPU and is a hardware problem.

The size of the MINIX process memory image increased by 10% (about 12K bytes). (Part of the increase, about 13%, is a result of increasing the number of process table slots from 16 to 24 to allow more processes to be active.)

Response time for interactive processes suffers when CPU bound processes are running. It becomes more prevalent with swapping because more processes can be run concurrently. The response degradation occurs because MINIX assigns the same priority to all user processes. A new process scheduler algorithm is required to remedy this.

5.2 Improvements and Further Development

The implementation provides an opportunity to devise, experiment, and test various scheduling and swapping algorithms. In addition, the implementation can be extended in other ways:

- User process priorities can be established to assign interactive processes a higher user priority to provide immediate response.
- Better swap out/in algorithms can be implemented, possibly integrating process scheduling with process swapping.

- Shared text can be implemented to make swapping more efficient.
- A real *brk()* system call can now be implemented in MINIX.
- The exec swap out can be optimized to eliminate a read and write of the incoming program.
- A more efficient swap device can be implemented.

OS	COMPILE TASK	HARDWARE	
		AT	XT
MINIX 1.3 v	kernel	7:43	
	FS	5:43	
	MM	3:22	
	TOTAL	16:48	
MINIX 1.3v w/swapping not used	kernel	7:45	49:09
	FS	5:53	
	MM	3:24	20:02
	TOTAL	17:02	69:11
MINIX 1.3 v w/swapping in use	kernel & MM	13:44	75:42
	kernel, MM, & FS	25:38	

TABLE 5-1. MINIX NON-SWAPPING vs SWAPPING BENCHMARKS

BIBLIOGRAPHY

- [1] Tanenbaum, A. S., "Operating Systems, Design and Implementation", Prentice-Hall, Inc., 1987
- [2] Bach, M. J., "The Design of the UNIX Operating System", Prentice-Hall, Inc., 1986
- [3] Lauer, H.C., Needham, R.M., "On the Duality of Operating System Structures" in Proc. Second International Symposium on Operating Systems, IRIA, Oct 1978, reprinted in Operating Systems Review, 13,2, April 1979, pp 3-19
- [4] Peachey, D. R., Bunt, R. B., Williamson, L.L., Brecht, T. B., "An Experimental Investigation on Scheduling Strategies for UNIX", Performance Evaluation Review, V12, 3 (August 1984), pp 158-166
- [5] Abdallah, M. S., "An Investigation of the Swapping Process in the UNIX Operating System", MSc. Thesis, Department of Computational Science, University of Saskatchewan, July 1982.
- [6] Peachey, R. B., Williamson, L.L., Bunt, R. B., "Taming the UNIX Scheduler", Proceedings CMG XV, San Francisco, Dec 1984

APPENDIX A - MODIFICATIONS TO KERNEL CODE

Appendix A-2 - MODIFICATIONS TO KERNEL CODE

Jul 6 14:30 1989 KERNEL.H Page 1

```

1          const.h
2
> 3      /* swap_task defines */
> 4      #define SWAP_IDLE 1      /* nothing to swap in or out */
> 5      #define SWAP_IN 2      /* runnable process on swap device */
> 6      #define SWAP_OUT 3      /* swap_task has ordered a swap out */
> 7      #define NOMES 4      /* No message is waiting for MM */
> 8      #define MPENDING 5      /* Message is waiting for MM */
9
10
11          glo.h
12
> 13      /* trace display */
> 14      EXTERN int Dflag;      /* if == 0, no display, else auto display */
> 15      EXTERN int swap_stat;
> 16      EXTERN struct sw_mm_mes swap_proc;
17
18
19          proc.h
20
> 21      unsigned res_time;      /* residence time in seconds (core, swap) */
> 22      #define SWAPPED 040      /* process is on swap device */
> 23      #define BLOCKED 0100      /* swapped proc is blocked */
> 24      #define STICKY 0200      /* proc has sticky bit set */
25
26
27          type.h
28
> 29      PUBLIC struct sw_mm_mes {
> 30          int status;
> 31          message *ms;
> 32      };

```

Appendix A-3 - MODIFICATIONS TO KERNEL CODE

Jul 6 15:13 1989 CLOCK.C Page 1

```

1  /*=====
2  *                                     do_clocktick                                     *
3  *=====*/
4  PRIVATE do_clocktick()
5  {
6  /* This routine is called on every clock tick. */
7
8  register struct proc *rp;
9  register int t, proc_nr;
10 extern int pr_busy, pcount, cum_count, prev_ct;
11
12 /* To guard against race conditions, first copy 'lost_ticks' to a local
13  * variable, add this to 'realtime', and then subtract it from 'lost_ticks'.
14  */
15 t = lost_ticks; /* 'lost_ticks' counts missed interrupts */
16 realtime += t + 1; /* update the time of day */
17 lost_ticks -= t; /* these interrupts are no longer missed */
18
19 if (next_alarm <= realtime) {
20 /* An alarm may have gone off, but proc may have exited, so check. */
21 next_alarm = MAX_P_LONG; /* start computing next alarm */
22 for (rp = &proc[0]; rp < &proc[NR_TASKS+NR_PROCS]; rp++) {
23 if (rp->p_alarm != (realtime) 0) {
24 /* See if this alarm time has been reached. */
25 if (rp->p_alarm <= realtime) {
26 /* A timer has gone off. If it is a user proc,
27  * send it a signal. If it is a task, call the
28  * function previously specified by the task.
29  */
30 proc_nr = rp - proc - NR_TASKS;
31 if (proc_nr >= 0)
32 cause_sig(proc_nr, SIGALRM);
33 else
34 (*watch_dog[proc_nr])();
35 rp->p_alarm = 0;
36 }
37
38 /* Work on determining which alarm is next. */
39 if (rp->p_alarm != 0 && rp->p_alarm < next_alarm)
40 next_alarm = rp->p_alarm;
41 }
42 }
43 }
44
45 accounting(); /* keep track of who is using the cpu */
46
47 /* If input characters are accumulating on an RS232 line, process them. */
48 if (flush_flag) {
49 t = (int) realtime; /* only low-order bits matter */
50 if ((t & FLUSH_MASK) == 0) rs_flush(); /* flush tty input */
51 }
52
53 /* If a user process has been running too long, pick another one. */

```

Appendix A-4 - MODIFICATIONS TO KERNEL CODE

Jul 6 15:13 1989 CLOCK.C Page 2

```

54     if (--sched_ticks == 0) {
55         if (bill_ptr == prev_ptr) sched(); /* process has run too long */
56         sched_ticks = SCHED_RATE;          /* reset quantum */
57         prev_ptr = bill_ptr;                /* new previous process */
58
59         /* Check if printer is hung up, and if so, restart it */
60         if (pr_busy && pcount > 0 && cum_count == prev_ct) pr_char();
61         prev_ct = cum_count; /* record # characters printed so far */
62     }
63
64     /* if residence update time, then update all residence times */
65     if (resd_ticks <= realtime) {
66         resd_ticks = realtime + RES_RATE;
67         for (rp = proc_addr(LOW_USER); rp < &proc[NR_TASKS+NR_PROCS]; rp++)
68             if (rp->p_flags != P_SLOT_FREE)
69                 rp->res_time++;
70         if (swap_stat == SWAP_IN) {
71             /* notify swap task that residence times have changed */
72             mess.m_source = CLOCK;
73             mess.m_type = CORE_IS_FREE;
74             send(SWAP_TASK, &mess);
75         }
76         if (seconds++ >= 5) {
77             seconds = 0;
78             if (Dflag)
79                 a_dmp();
80         }
81     }
82 }

```

Appendix A-5 - MODIFICATIONS TO KERNEL CODE

Jul 6 15:13 1989 PROC.C Page 1

```

1  /*=====
2  *                               mini_rec                               *
3  *=====*/
4  PRIVATE int mini_rec(caller, src, m_ptr)
5  int caller; /* process trying to get message */
6  int src; /* which message source is wanted (or ANY) */
7  message *m_ptr; /* pointer to message buffer */
8  {
9  /* A process or task wants to get a message. If one is already queued,
10 * acquire it and deblock the sender. If no message from the desired source
11 * is available, block the caller. No need to check parameters for validity.
12 * Users calls are always sendrec(), and mini_send() has checked already.
13 * Calls from the tasks, MM, and FS are trusted.
14 */
15
16 extern struct sw_mm_mes swap_proc;
17 register struct proc *caller_ptr, *sender_ptr, *previous_ptr;
18 int sender;
19
20 caller_ptr = proc_addr(caller); /* pointer to caller's proc struct */
21
22 /* Check to see if a message from desired source is already available. */
23 sender_ptr = caller_ptr->p_callerq;
24 if ((caller_ptr->p_flags & SENDING) == 0) {
25 while (sender_ptr != NIL_PROC) {
26 sender = sender_ptr - proc - NR_TASKS;
27 if (src == ANY || src == sender) {
28 /* An acceptable message has been found. */
29 cp_mess(sender, sender_ptr->p_map[D].mem_phys,
30 sender_ptr->p_messbuf, caller_ptr->p_map[D].mem_phys, m_ptr);
31 sender_ptr->p_flags &= SENDING; /* deblock sender */
32 if (sender_ptr->p_flags == 0) ready(sender_ptr);
33 if (sender_ptr == caller_ptr->p_callerq)
34 caller_ptr->p_callerq = sender_ptr->p_sendlink;
35 else
36 previous_ptr->p_sendlink = sender_ptr->p_sendlink;
37 return(OK);
38 }
39 previous_ptr = sender_ptr;
40 sender_ptr = sender_ptr->p_sendlink;
41 }
42 }
43
44 /* No suitable message is available. Block the process trying to receive. */
45 caller_ptr->p_getfrom = src;
46 caller_ptr->p_messbuf = m_ptr;
47 if (caller_ptr->p_flags == 0) unready(caller_ptr);
48 caller_ptr->p_flags |= RECEIVING;
49
50 /* If MM has just blocked and there are kernel signals pending, now is the
51 * time to tell MM about them, since it will be able to accept the message.
52 */
53 if (((sig_procs > 0) || (swap_proc.status == MPENDING)) &&

```

Appendix A-6 - MODIFICATIONS TO KERNEL CODE

Jul 6 15:13 1989 PROC.C Page 2

```
> 54      (caller == MM_PROC_NR && src == ANY) ) {  
55          inform();  
56          pick_proc();  
57      }  
58      return(OK);  
59  }
```

Appendix A-7 - MODIFICATIONS TO KERNEL CODE

Jul 6 15:13 1989 SWAPPER.C Page 1

```

> 1  /* This file contains the code and data for the swapper task. It
> 2  * has a single entry point, swap_task(). It accepts six message
> 3  * types:
> 4  *
> 5  * CORE_IS_FREE: memory has been released or can possibly be freed
> 6  * CORE_IS_NEEDED: runnable process is on swap device
> 7  * SWAP_IN_COMPL: swap_task request has succeeded
> 8  * SWAP_IN_FAILED: swap_task request has failed
> 9  * SWAP_OUT_COMPL: swap_task or MM request has succeeded
> 10 * SWAP_OUT_FAILED: swap_task request has failed
> 11 *
> 12 * The input message is format m1. The parameters are:
> 13 *
> 14 * m_type PROC1 PROC2 PID MEM_PTR UTILITY
> 15 *
> 16 * CORE_IS_FREE | | | | | |
> 17 * |-----+-----+-----+-----+-----+-----|
> 18 * CORE_IS_NEEDED |unblocked| | | | |
> 19 * |-----+-----+-----+-----+-----+-----|
> 20 * SWAP_IN_COMPL |swappedin| | | | |
> 21 * |-----+-----+-----+-----+-----+-----|
> 22 * SWAP_OUT_COMPL |swapedout| blocked | | | |
> 23 * |-----+-----+-----+-----+-----+-----|
> 24 * SWAP_OUT_FAILED | proc no | | | | |
> 25 * |-----+-----+-----+-----+-----+-----|
> 26 *
> 27 * The input message is format m2. The parameters are:
> 28 *
> 29 * | SWAP_IN_FAILED |sizeofcor|pausewait| | | |
> 30 * |-----+-----+-----+-----+-----+-----|
> 31 */
> 32
> 33 #include "../const.h"
> 34 #include "../type.h"
> 35 #include "../callnr.h"
> 36 #include "../com.h"
> 37 #include "../error.h"
> 38 #include "../signal.h"
> 39 #include "const.h"
> 40 #include "type.h"
> 41 #include "glo.h"
> 42 #include "proc.h"
> 43
> 44 /* constant definitions */
> 45 #define NOP 1 /* nothing in progress */
> 46 #define SW_INP 2 /* swap-in in progress */
> 47 #define SW_OUTP 3 /* swap-out in progress */
> 48
> 49 /* swap out algorithm definitions */
> 50 #define PICKSIZE 10000 /* size is >= size needed */
> 51 #define PICKSTICK 1000 /* sticky bit is not set */
> 52 #define PICKPWBK 100 /* proc is PAUSE/WAIT */
> 53 #define PICKSRBLK 10 /* proc is blocked on both S & R */

```

Appendix A-8 - MODIFICATIONS TO KERNEL CODE

Jul 6 15:13 1989 SWAPPER.C Page 2

```

> 54 #define PICKNBLK 1 /* proc is not blocked at all */
> 55 #define PICKINELL -20000 /* proc is ineligible for swap out */
> 56
> 57 #define CRESMIN 15 /* min residence time in seconds */
> 58 #define LONG1 m2_11 /* message slot to carry long */
> 59
> 60 /* swapper task variables */
> 61 PRIVATE int inprogres;
> 62 PRIVATE message mes, muns;
> 63
> 64
> 65 /*=====
> 66 * swap_task *
> 67 *=====*/
> 68 PUBLIC swap_task()
> 69 {
> 70 /* Main program of swap task. It determines which of the 6 possible
> 71 * calls this is by looking at 'mes.m_type'. Then it dispatches.
> 72 */
> 73 struct proc *rp;
> 74 int opcode, pnum;
> 75 long brn;
> 76 phys_clicks sizeneed;
> 77 phys_bytes src_phys, dst_phys;
> 78 vir_bytes ptr;
> 79
> 80 init_swap(); /* initialize swap tables */
> 81
> 82 /* Main loop of the swap task. Get work, process it, sometimes reply. */
> 83 while (TRUE) {
> 84 receive(ANY, &mes); /* go get a message */
> 85 opcode = mes.m_type; /* extract the message type */
> 86
> 87 switch (opcode) {
> 88 case CORE_IS_FREE:
> 89 if (swap_stat == SWAP_IN) && (inprogres == NOP) {
> 90 try_to_swin();
> 91 }
> 92 break;
> 93
> 94 case CORE_IS_NEEDED:
> 95 if (mes.PROC1 != 0) {
> 96 /* PROC1 on swap device just became unblocked */
> 97 rp = proc_addr(mes.PROC1);
> 98 rp->p_flags &= BLOCKED;
> 99 }
> 100 if (swap_stat == SWAP_IDLE) {
> 101 try_to_swin();
> 102 }
> 103 break;
> 104
> 105 case SWAP_IN_COMPL:
> 106 if (inprogres == SW_INP) {

```

Appendix A-9 - MODIFICATIONS TO KERNEL CODE

Jul 6 15:13 1989 SWAPPER.C Page 3

```

> 107             inprogres = NOP;
> 108             pnun = mes.PROC1;
> 109             rp = proc_addr(pnun);
> 110             rp->p_flags &= SWAPPED;
> 111             rp->res_time = 0;
> 112             if (rp->p_flags == 0) ready(rp);
> 113             try_to_swin0;
> 114         } else
> 115             printf("S_I_C: swpin NOT in progress0);
> 116         break;
> 117
> 118     case SWAP_IN_FAILED:
> 119         if(inprogres == SW_INP) {
> 120             inprogres = NOP;
> 121             sizeneed = (phys_clicks)mes.PROC1;
> 122             bm = mes.LONG1;
> 123             mms.m_source = SWAP_TASK;
> 124             if( (pnun = pik_outsw(bm, sizeneed)) != 0) {
> 125                 rp = proc_addr(pnun);
> 126                 rp->p_flags |= SWAPPED;
> 127                 unready(rp);
> 128                 inprogres = SW_OUTP;
> 129                 mms.m_type = SWAP_OUT_REQ;
> 130                 mms.PROC1 = pnun;
> 131             } else {
> 132                 mms.m_type = NO_SWAP_OUT;
> 133             }
> 134             send(MM_PROC_NR, &mms);
> 135         } else
> 136             printf("S_I_F: swap in NOT in progress0);
> 137         break;
> 138
> 139     case SWAP_OUT_COMPL:
> 140         if(inprogres == SW_OUTP) {
> 141             inprogres = NOP;
> 142             rp = proc_addr(mes.PROC1);
> 143             rp->p_flags |= NO_MAP;
> 144             if(mes.PROC2)
> 145                 rp->p_flags |= BLOCKED;
> 146             rp->res_time = 0;
> 147             try_to_swin0;
> 148         } else
> 149             printf("S_O_C: swap out NOT in progress0);
> 150         break;
> 151
> 152     case SWAP_OUT_FAILED:
> 153         if(inprogres == SW_OUTP) {
> 154             rp = proc_addr(mes.PROC1);
> 155             rp->p_flags &= SWAPPED;
> 156             inprogres = NOP;
> 157             if (rp->p_flags == 0) ready(rp);
> 158         } else
> 159             printf("S_O_F: swap out NOT in progress0);

```


Appendix A-10 - MODIFICATIONS TO KERNEL CODE

Jul 6 15:13 1989 SWAPPER.C Page 4

```

> 160         break;
> 161
> 162         default: panic("swap task got bad message", mes.m_type);
> 163         break;
> 164     }
> 165 }
> 166 }
> 167
> 168 init_swap()
> 169 {
> 170     inprogres = NOP;
> 171     swap_proc.status = NOMES;
> 172     swap_stat = SWAP_IDLE;
> 173 }
> 174
> 175 /* set up non-blocking message transfer to MM */
> 176 set_swapproc(type, num)
> 177 int type, num;
> 178 {
> 179     mms.m_source = SWAP_TASK;
> 180     mms.m_type = type;
> 181     mms.PROC1 = num;
> 182     swap_proc.ms = &mms;
> 183     swap_proc.status = MPENDING;
> 184     if( ((proc[NR_TASKS + MM_PROC_NR].p_flags & RECEIVING) == 0) ||
> 185         (proc[NR_TASKS + MM_PROC_NR].p_getfrom != ANY) )
> 186         return;
> 187     inform();
> 188 }
> 189
> 190 try_to_swin()
> 191 {
> 192     int proc;
> 193     if( (proc = pik_insw()) == 0 ) {
> 194         swap_stat = SWAP_IDLE;
> 195     } else {
> 196         inprogres = SW_INP;
> 197         swap_stat = SWAP_IN;
> 198         set_swapproc(SWAP_IN_REQ, proc);
> 199     }
> 200 }
> 201
> 202 pik_insw()
> 203 {
> 204     /* returns the proc slot # of the most eligible proc on swap device */
> 205     /* to swap in. If none are eligible, then 0 is returned */
> 206     struct proc *rp;
> 207     struct proc *stkrp;
> 208     struct proc *nostkrp;
> 209     unsigned stkres, nostkres;
> 210     int pick;
> 211
> 212     stkrp = proc_addr(0);

```

Appendix A-11 - MODIFICATIONS TO KERNEL CODE

Jul 6 15:13 1989 SWAPPER.C Page 5

```

> 213     nostkrp = proc_addr(0);
> 214     stkrp = 0;
> 215     nostkres = 0;
> 216
> 217     for(rp = proc_addr(INIT_PROC_NR + 1); rp < proc_addr(NR_PROCS); rp++) {
> 218         if(rp->p_flags & P_SLOT_FREE) continue;
> 219         if((rp->p_flags & SWAPPED) &&
> 220             ((rp->p_flags & BLOCKED) == 0) ) {
> 221             if(rp->p_flags & STICKY) {
> 222                 if(rp->res_time >= stkrp) {
> 223                     stkrp = rp;
> 224                     stkrp = rp->res_time;
> 225                 }
> 226             } else if(rp->res_time >= nostkres) {
> 227                 nostkrp = rp;
> 228                 nostkres = rp->res_time;
> 229             }
> 230         }
> 231     }
> 232     if(stkrp != proc_addr(0))
> 233         pick = (int)(stkrp - proc_addr(0));
> 234     else
> 235         pick = (int)(nostkrp - proc_addr(0));
> 236     return(pick);
> 237 }
> 238
> 239
> 240     pik_outsw(map, sizeneed)
> 241     long map;
> 242     int sizeneed;
> 243     {
> 244         /* returns the proc slot # of the most eligible proc in core */
> 245         /* to swap out. If none are eligible, then 0 is returned */
> 246
> 247         struct proc *rp, *hrp;
> 248         int priority, hpriority;
> 249         unsigned hres;
> 250
> 251         hpriority = 0;
> 252         hres = 0;
> 253         hrp = proc_addr(0);
> 254         for(rp = proc_addr(INIT_PROC_NR + 1); rp < proc_addr(NR_PROCS); rp++) {
> 255             if(rp->p_flags & (P_SLOT_FREE | SWAPPED)) continue;
> 256             priority = 0;
> 257             if((rp->p_map[S].mem_phys +
> 258                 rp->p_map[S].mem_len -
> 259                 rp->p_map[T].mem_phys) >= sizeneed)
> 260                 priority += PICKSIZE;
> 261             if((rp->p_flags & STICKY) == 0)
> 262                 priority += PICKSTICK;
> 263             priority += blocktype(rp, map);
> 264             if( (priority > hpriority) ||
> 265                 (priority == hpriority && rp->res_time > hres) ) {

```

Appendix A-12 - MODIFICATIONS TO KERNEL CODE

Jul 6 15:13 1989 SWAPPER.C Page 6

```

> 266             hpriority = priority;
> 267             hres = rp->res_time;
> 268             hrp = rp;
> 269         }
> 270     }
> 271     if(hpriority == 0)
> 272         return(0);
> 273     return( (int)(hrp - proc_addr(0)));
> 274 }
> 275
> 276
> 277 /* determine if process is blocked or pause/wait, SEND & RECEIVE, */
> 278 /* or not at all */
> 279 hlocktype(rp, map)
> 280 struct proc *rp;
> 281 long map;
> 282 {
> 283
> 284     if( (map >> (rp - proc_addr(0)) & 1) &&
> 285         (rp->res_time >= CRESMIN) )
> 286         return(PICKPWBK);
> 287     if( (rp->p_flags & (SENDING | RECEIVING) == 0) &&
> 288         (rp->res_time >= CRESMIN) )
> 289         return(PICKNBK);
> 290     if( (rp->p_flags & (SENDING | RECEIVING) == (SENDING | RECEIVING)) &&
> 291         (rp->res_time >= CRESMIN) )
> 292         return(PICKINELL);
> 293     return(PICKINELL);
> 294 }

```

Appendix A-13 - MODIFICATIONS TO KERNEL CODE

Jul 6 15:13 1989 SYSTEM.C Page 1

```

1  /*=====
2  *                                     *
3  *                                     * do_fork
4  *=====*/
5  PRIVATE int do_fork(m_ptr)
6  message *m_ptr; /* pointer to request message */
7  {
8  /* Handle sys_fork(). 'k1' has forked. The child is 'k2'. */
9
10 register struct proc *rpc;
11 register char *sptr, *dptr; /* pointers for copying proc struct */
12 int k1; /* number of parent process */
13 int k2; /* number of child process */
14 int pid; /* process id of child */
15 int bytes; /* counter for copying proc struct */
16 int fkswap; /* TRUE, if fork has swapped out child */
17
18 k1 = m_ptr->PROC1; /* extract parent slot number from msg */
19 k2 = m_ptr->PROC2; /* extract child slot number */
20 pid = m_ptr->PID; /* extract child process id */
21 fkswap = (int)m_ptr->UTILITY; /* extract swap status of child */
22
23 if (k1 < 0 || k1 >= NR_PROCS || k2 < 0 || k2 >= NR_PROCS) return(E_BAD_PROC);
24 rpc = proc_addr(k2);
25
26 /* Copy parent 'proc' struct to child. */
27 sptr = (char *) proc_addr(k1); /* parent pointer */
28 dptr = (char *) proc_addr(k2); /* child pointer */
29 bytes = sizeof(struct proc); /* # bytes to copy */
30 while (bytes-- > 0) *dptr++ = *sptr++; /* copy parent struct to child */
31
32 rpc->p_flags |= NO_MAP; /* inhibit the process from running */
33 rpc->p_flags &= PENDING; /* only one in group should have PENDING */
34 rpc->p_pending = 0;
35 rpc->p_pid = pid; /* install child's pid */
36 rpc->p_reg[RET_REG] = 0; /* child sees pid = 0 to know it is child */
37
38 rpc->user_time = 0; /* set all the accounting times to 0 */
39 rpc->sys_time = 0;
40 rpc->res_time = 0;
41 rpc->child_uptime = 0;
42 rpc->child_stime = 0;
43 if(fkswap) {
44     rpc->p_flags |= SWAPPED;
45     if(swap_stat == SWAP_IDLE) {
46         mess.m_source = SYSTASK;
47         mess.m_type = CORE_IS_NEEDED;
48         mess.PROC1 = 0;
49         send(SWAP_TASK, &mess);
50     }
51 }
52 return(OK);
53 }

```

Appendix A-14 - MODIFICATIONS TO KERNEL CODE

Jul 6 15:13 1989 SYSTEM.C Page 2

```

54
55  /* =====
56  *                               do_newmap                               *
57  * ===== */
58 PRIVATE int do_newmap(m_ptr)
59 message *m_ptr;                /* pointer to request message */
60 {
61     /* Handle sys_newmap(). Fetch the memory map from MM. */
62
63     register struct proc *rp, *rsr;
64     phys_bytes src_phys, dst_phys, pn;
65     vir_bytes vmm, vsys, vn;
66     int caller;                /* whose space has the new map (usually MM) */
67     int k;                     /* process whose map is to be loaded */
68     int old_flags;             /* value of flags before modification */
69     struct mem_map *map_ptr;   /* virtual address of map inside caller (MM) */
70     int util;                  /* TRUE, if exec uses core as a buffer */
71
72     /* Extract message parameters and copy new memory map from MM. */
73     caller = m_ptr->m_source;
74     k = m_ptr->PROC1;
75     map_ptr = (struct mem_map *) m_ptr->MEM_PTR;
76     util = (int)m_ptr->UTILITY; /* extract swap status of child */
77     if (k < NR_TASKS || k >= NR_PROCS) return(E_BAD_PROC);
78     rp = proc_addr(k);        /* ptr to entry of user getting new map */
79     rsr = proc_addr(caller);   /* ptr to MM's proc entry */
80     vn = NR_SEGS * sizeof(struct mem_map);
81     pn = vn;
82     vmm = (vir_bytes) map_ptr; /* careful about sign extension */
83     vsys = (vir_bytes) rp->p_map; /* again, careful about sign extension */
84     if ((src_phys = umap(rsr, D, vmm, vn)) == 0)
85         panic("bad call to sys_newmap (src)", NO_NUM);
86     if ((dst_phys = umap(proc_addr(SYSTASK), D, vsys, vn)) == 0)
87         panic("bad call to sys_newmap (dst)", NO_NUM);
88     phys_copy(src_phys, dst_phys, pn);
89
90     #ifdef i8088
91     /* On 8088, set segment registers. */
92     rp->p_reg[CS_REG] = rp->p_map[T].mem_phys; /* set cs */
93     rp->p_reg[DS_REG] = rp->p_map[D].mem_phys; /* set ds */
94     rp->p_reg[SS_REG] = rp->p_map[D].mem_phys; /* set ss */
95     rp->p_reg[ES_REG] = rp->p_map[D].mem_phys; /* set es */
96     #endif
97
98     /* don't make process ready if core is being used as an I/O buffer */
99     if (!util) {
100         old_flags = rp->p_flags; /* save the previous value of the flags */
101         rp->p_flags &= NO_MAP;
102         if (old_flags != 0 && rp->p_flags == 0) ready(rp);
103     } else {
104         rp->p_flags |= NO_MAP;
105     }
106     return(OK);

```

Appendix A-15 - MODIFICATIONS TO KERNEL CODE

Jul 6 15:13 1989 SYSTEM.C Page 3

```

107 }
108
109
110 /*=====
111  *                               do_exec                               *
112  *=====*/
113 PRIVATE int do_exec(m_ptr)
114 message *m_ptr; /* pointer to request message */
115 {
116 /* Handle sys_exec(). A process has done a successful EXEC. Patch it up. */
117
118 register struct proc *rp;
119 int k; /* which process */
120 int *sp; /* new sp */
121 int exswap; /* TRUE, if exec has swapped out new proc */
122
123 k = m_ptr->PROC1; /* 'k' tells which process did EXEC */
124 sp = (int *) m_ptr->STACK_PTR;
125 exswap = (int)m_ptr->UTILITY; /* extract swap status of new proc */
126 if (k < 0 || k >= NR_PROCS) return(E_BAD_PROC);
127 rp = proc_addr(k);
128 rp->p_sp = sp; /* set the stack pointer */
129 rp->p_cpsw.pc = (int (*)0) 0; /* reset pc */
130 rp->p_alarm = 0; /* reset alarm timer */
131 rp->p_flags &= RECEIVING; /* MM does not reply to EXEC call */
132
133 if(exswap) {
134     rp->p_flags |= SWAPPED;
135     if(swap_stat == SWAP_IDLE) {
136         mess.m_source = SYSTASK;
137         mess.m_type = CORE_IS_NEEDED;
138         mess.PROC1 = 0;
139         send(SWAP_TASK, &mess);
140     }
141     } else if(swap_stat == SWAP_IN) {
142         /* notify swapper that core has been freed */
143         mess.m_source = SYSTASK;
144         mess.m_type = CORE_IS_FREE;
145         send(SWAP_TASK, &mess);
146     }
147
148 if (rp->p_flags == 0) ready(rp);
149 if(exswap)
150     set_name(k, (char *)0); /* erase command string from F1 display */
151 else
152     set_name(k, (char *)sp); /* save command string for F1 display */
153
154 return(OK);
155 }
156
157
158 /*=====
159  *                               do_xit                               *
160  *=====*/

```

Appendix A-16 - MODIFICATIONS TO KERNEL CODE

Jul 6 15:13 1989 SYSTEM.C Page 4

```

160  *-----*/
161  PRIVATE int do_xit(m_ptr)
162  message *m_ptr; /* pointer to request message */
163  {
164  /* Handle sys_xit(). A process has exited. */
165
166  register struct proc *rp, *rc;
167  struct proc *np, *xp;
168  int parent; /* number of exiting proc's parent */
169  int proc_nr; /* number of process doing the exit */
170
171  parent = m_ptr->PROC1; /* slot number of parent process */
172  proc_nr = m_ptr->PROC2; /* slot number of exiting process */
173  if (parent < 0 || parent >= NR_PROCS || proc_nr < 0 || proc_nr >= NR_PROCS)
174      return(E_BAD_PROC);
175  rp = proc_addr(parent);
176  rc = proc_addr(proc_nr);
177  rp->child_etime += rc->user_time + rc->child_etime; /* accum child times */
178  rp->child_stime += rc->sys_time + rc->child_stime;
179  rc->p_alarm = 0; /* turn off alarm timer */
180  if (rc->p_flags == 0) unready(rc);
181  set_name(proc_nr, (char *) 0); /* disable command printing for FI */
182
183  /* If the process being terminated happens to be queued trying to send a
184  * message (i.e., the process was killed by a signal, rather than it doing an
185  * EXIT), then it must be removed from the message queues.
186  */
187  if (rc->p_flags & SENDING) {
188      /* Check all proc slots to see if the exiting process is queued. */
189      for (rp = &proc[0]; rp < &proc[NR_TASKS + NR_PROCS]; rp++) {
190          if (rp->p_callerq == NIL_PROC) continue;
191          if (rp->p_callerq == rc) {
192              /* Exiting process is on front of this queue. */
193              rp->p_callerq = rc->p_sendlink;
194              break;
195          } else {
196              /* See if exiting process is in middle of queue. */
197              np = rp->p_callerq;
198              while ( ( xp = np->p_sendlink ) != NIL_PROC )
199                  if (xp == rc) {
200                      np->p_sendlink = xp->p_sendlink;
201                      break;
202                  } else {
203                      np = xp;
204                  }
205          }
206      }
207  }
208  if (rc->p_flags & PENDING) --sig_procs;
209  rc->p_flags = P_SLOT_FREE;
210
211  /* notify swapper that core has been freed */
212  if (swap_stat == SWAP_IN) {

```

Appendix A-17 - MODIFICATIONS TO KERNEL CODE

Jul 6 15:13 1989 SYSTEM.C Page 5

```

> 213     mess.m_type = CORE_IS_FREE;
> 214     send(SWAP_TASK, &mess);
> 215 }
> 216
> 217     return(OK);
> 218 }
> 219
> 220
> 221 /*=====
> 222 *                               do_lock                               *
> 223 *=====*/
> 224 PRIVATE int do_lock(m_ptr)
> 225     message *m_ptr;           /* pointer to request message */
> 226 {
> 227     /* Handle sys_lock or restore request from MM */
> 228
> 229     if(m_ptr->LOCK_RES == LOCK) {
> 230         m_ptr->PSW = (int) lock();
> 231     } else {
> 232         restore( (unsigned) m_ptr->PSW);
> 233     }
> 234     return(OK);
> 235 }
> 236
> 237
> 238 /*=====
> 239 *                               inform                               *
> 240 *=====*/
> 241 PUBLIC inform()
> 242 {
> 243     /* When a signal is detected by the kernel (e.g., DEL), or generated by a task
> 244     * (e.g. clock task for SIGALRM), cause_sig() is called to set a bit in the
> 245     * p_pending field of the process to signal. Then inform() is called to see
> 246     * if MM is idle and can be told about it. Whenever MM blocks, a check is
> 247     * made to see if 'sig_procs' is nonzero; if so, inform() is called.
> 248     */
> 249
> 250     register struct proc *rp;
> 251
> 252     /* MM is waiting for new input. Find a process with pending signals. */
> 253     /* does swapper want to send message to MM now? */
> 254     if(swap_proc.status == MPENDING) {
> 255         if (mini_send(SWAP_TASK, MM_PROC_NR, swap_proc.ms) != OK)
> 256             panic("can't inform MM", NO_NUM);
> 257         swap_proc.status = NOMES;
> 258         return;
> 259     }
> 260     for (rp = proc_addr(0); rp < proc_addr(NR_PROCS); rp++)
> 261         if (rp->p_flags & PENDING) {
> 262             m.m_type = KSIG;
> 263             m.PROC1 = rp - proc - NR_TASKS;
> 264             m.SIG_MAP = rp->p_pending;
> 265             sig_procs--;

```


Appendix A-18 - MODIFICATIONS TO KERNEL CODE

Jul 6 15:13 1989 SYSTEM.C Page 6

```
266             if (mini_send(HARDWARE, MM_PROC_NR, &m) != OK)
267                 panic("can't inform MM", NO_NUM);
268             rp->p_pending = 0; /* the ball is now in MM's court */
269             rp->p_flags &= PENDING;
270             if (rp->p_flag == 0) ready(rp);
271             return;
272         }
273     }
```

Appendix A-19 - MODIFICATIONS TO KERNEL CODE

Jul 6 15:13 1989 TABLE.C Page 1

```

1  /* The object file of "table.c" contains all the data. In the *.h files,
2  * declared variables appear with EXTERN in front of them, as in
3  *
4  *   EXTERN int x;
5  *
6  * Normally EXTERN is defined as extern, so when they are included in another
7  * file, no storage is allocated. If the EXTERN were not present, but just
8  * say,
9  *
10 *   int x;
11 *
12 * then including this file in several source files would cause 'x' to be
13 * declared several times. While some linkers accept this, others do not,
14 * so they are declared extern when included normally. However, it must
15 * be declared for real somewhere. That is done here, by redefining
16 * EXTERN as the null string, so the inclusion of all the *.h files in
17 * table.c actually generates storage for them. All the initialized
18 * variables are also declared here, since
19 *
20 *   extern int x = 4;
21 *
22 * is not allowed. If such variables are shared, they must also be declared
23 * in one of the *.h files without the initialization.
24 */
25
26 #include "../h/const.h"
27 #include "../h/type.h"
28 #include "../h/com.h"
29 #include "const.h"
30 #include "type.h"
31 #undef EXTERN
32 #define EXTERN
33 #include "glo.h"
34 #include "proc.h"
35 #include "tty.h"
36
37 extern int sys_task(), clock_task(), mem_task(), floppy_task(),
> 38 winchester_task(), tty_task(), printer_task(), swap_task();
39
40 #ifdef AM_KERNEL
41 extern int amoeba_task();
42 extern int amint_task();
43 #endif
44
45 /* The startup routine of each task is given below, from -NR_TASKS upwards.
46 * The order of the names here MUST agree with the numerical values assigned to
47 * the tasks in ../h/com.h.
48 */
> 49 #define SMALL_STACK 512
50
51 #define TTY_STACK SMALL_STACK
52
> 53 #define SWAP_STACK SMALL_STACK

```

Appendix A-20 - MODIFICATIONS TO KERNEL CODE

Jul 6 15:13 1989 TABLE.C Page 2

```

54
55 #define PRINTER_STACK SMALL_STACK
56 #define WINCH_STACK SMALL_STACK
57 #define FLOP_STACK SMALL_STACK
58 #define MEM_STACK SMALL_STACK
59 #define CLOCK_STACK SMALL_STACK
60 #define SYS_STACK SMALL_STACK
61
62
63
64 #ifdef AM_KERNEL
65 #   define AMINT_STACK SMALL_STACK
66 #   define AMOEBA_STACK 1532
67 #   define AMOEBA_STACK_SPACE (AM_NTASKS*AMOEBA_STACK + AMINT_STACK)
68 #else
69 #   define AMOEBA_STACK_SPACE 0
70 #endif
71
72 #define TOT_STACK_SPACE (TTY_STACK + AMOEBA_STACK_SPACE + \
> 73 SWAP_STACK + \
74 PRINTER_STACK + \
75 WINCH_STACK + FLOP_STACK + \
76 MEM_STACK + CLOCK_STACK + SYS_STACK)
77
78 /*
79 ** some notes about the following table:
80 ** 1) The tty_task should always be first so that other tasks can use printf
81 ** if their initialisation has problems.
82 ** 2) If you add a new kernel task, add it after the amoeba_tasks and before
83 ** the printer task.
84 ** 3) The task name is used for process status (F1 key) and must be six (6)
85 ** characters in length. Pad it with blanks if it is too short.
86 */
87
88 PUBLIC struct tasktab tasktab[] = {
89     #ifdef AM_KERNEL
90         amint_task, AMINT_STACK, "AMINT ",
91         amoeba_task, AMOEBA_STACK, "AMTASK",
92         amoeba_task, AMOEBA_STACK, "AMTASK",
93         amoeba_task, AMOEBA_STACK, "AMTASK",
94         amoeba_task, AMOEBA_STACK, "AMTASK",
95     #endif
> 96     swap_task, SWAP_STACK, "SWAPER",
97     printer_task, PRINTER_STACK, "PRINTR",
98     winchester_task, WINCH_STACK, "WINCHE",
99     floppy_task, FLOP_STACK, "FLOPPY",
100     mem_task, MEM_STACK, "RAMDSK",
101     clock_task, CLOCK_STACK, "CLOCK ",
102     sys_task, SYS_STACK, "SYS ",
103     0, 0, "IDLE ",
104     0, 0, "MM ",
105     0, 0, "FS ",
106     0, 0, "INIT "

```

Appendix A-21 - MODIFICATIONS TO KERNEL CODE

Jul 6 15:13 1989 TABLE.C Page 3

```

107     };
108
109     int t_stack[TOT_STACK_SPACE/sizeof (int)];
110
111     int k_stack[K_STACK_BYTES/sizeof (int)];          /* The kernel stack. */
112
113
114     /*
115     ** The number of kernel tasks must be the same as NR_TASKS.
116     ** If NR_TASKS is not correct then you will get the compile error:
117     **   multiple case entry for value 0
118     ** The function ____dummy is never called.
119     */
120
121     #define NKT (sizeof tasktab / sizeof (struct tasktab) - (INIT_PROC_NR + 1))
122     ____dummy()
123     {
124         switch(0)
125         {
126             case 0:
127                 case (NR_TASKS == NKT):
128                     ;
129             }
130     }

```

APPENDIX B - MODIFICATIONS TO MEMORY MANAGER CODE

Appendix B-2 - MODIFICATIONS TO MEMORY MANAGER CODE

Jul 6 14:30 1989 MM.H Page 1

```
1          const.h
2
> 3  #define SWAP_MODE      0777      /* mode to use on swap device files */
4
5
6          mproc.h
7
> 8  unsigned mp_ssw_map;      /* bitmap of signals recvd while swapped */
> 9  int mp_deadchild;        /* >0 means WAITING SWAPPED proc's child died*/
>10  #define SWAPPED          0100     /* process is swapped out */
>11  #define FKSWAPPED        0200     /* process is swapped out by fork */
>12  #define WASPWS           0400     /* process was P/W & SWAPPED & then awakened */
```

Appendix B-3 - MODIFICATIONS TO MEMORY MANAGER CODE

Jul 6 14:41 1989 ALLOC.C Page 1

```

1  /*=====
2  *
3  * alloc_mem
4  *=====*/
5  PUBLIC phys_clicks alloc_mem(clicks)
6  phys_clicks clicks; /* amount of memory requested */
7  {
8  /* Allocate a block of memory from the free list using first fit. The block
9  * consists of a sequence of contiguous bytes, whose length in clicks is
10 * given by 'clicks'. A pointer to the block is returned. The block is
11 * always on a click boundary. This procedure is called when memory is
12 * needed for FORK or EXEC.
13 */
14 register struct hole *hp, *prev_ptr;
15 phys_clicks old_base;
16
17 while(TRUE) {
18     hp = hole_head;
19     while (hp != NIL_HOLE) {
20         if (hp->h_len >= clicks) {
21             /* We found a hole that is big enough. Use it. */
22             old_base = hp->h_base; /* remember where it started */
23             hp->h_base += clicks; /* bite a piece off */
24             hp->h_len -= clicks; /* ditto */
25
26             /* If hole is only partly used, reduce size and return. */
27             if (hp->h_len != 0) return(old_base);
28
29             /* The entire hole has been used up. Manipulate free list. */
30             del_slot(prev_ptr, hp);
31             return(old_base);
32         }
33
34         prev_ptr = hp;
35         hp = hp->h_next;
36     }
37     if (tot_hole() < clicks)
38         break;
39     compact(); /* mem is available, compact to get it */
40 }
41 return(NO_MEM);
42 }
43
44 /*=====
45 *
46 * tot_hole
47 *=====*/
48 PUBLIC phys_clicks tot_hole()
49 {
50 /* Scan the hole list and return sum of all holes. */
51
52 register struct hole *hp;
53 register phys_clicks total;

```

Appendix B-4 - MODIFICATIONS TO MEMORY MANAGER CODE

Jul 6 14:41 1989 ALLOC.C Page 2

```

> 54
> 55     hp = hole_head;
> 56     total = 0;
> 57     do {
> 58         total += hp->h_len;
> 59     }while ((hp = hp->h_next) != NIL_HOLE);
> 60     return(total);
> 61 }
> 62
> 63
> 64 /*=====
> 65      *                               *
> 66      *               compact               *
> 67      *                               *
> 68      *=====*/
> 69 PUBLIC compact()
> 70 {
> 71     /* Go through the memory hole map and eliminate all holes except one
> 72      * by actually moving process images into the empty spaces.
> 73      */
> 74     phys_clicks hole_base, old_base, proc_size;
> 75     long hb, ob, ps;
> 76     struct mproc *rmp;
> 77     unsigned old_state;
> 78     int found;
> 79
> 80     sys_lock(LOCK, &old_state);
> 81     while(hole_head->h_next != NIL_HOLE) {
> 82         /* find first hole */
> 83         hole_base = hole_head->h_base;
> 84
> 85         /* find proc just above this hole */
> 86         found = 0;
> 87         for(rmp = &mproc[INIT_PROC_NR + 1]; rmp < &mproc[NR_PROCS]; rmp++) {
> 88             if( ((rmp->mp_flags & IN_USE) == 0) ||
> 89                 (rmp->mp_flags & (SWAPPED | FKSWAPPED)) ||
> 90                 (rmp->mp_flags & HANGING) )
> 91                 continue;
> 92             if( (rmp->mp_seg[T].mem_phys - hole_base) == hole_head->h_len ) {
> 93                 found = 1;
> 94                 break;
> 95             }
> 96         }
> 97         if(found == 0) {
> 98             return;
> 99         }
> 100         old_base = rmp->mp_seg[T].mem_phys;
> 101         proc_size = rmp->mp_seg[S].mem_phys + rmp->mp_seg[S].mem_len - old_base;
> 102         hb = (long) hole_base << CLICK_SHIFT;
> 103         ob = (long) old_base << CLICK_SHIFT;
> 104         ps = (long) proc_size << CLICK_SHIFT;
> 105
> 106         if(mem_copy(ABS,0,ob, ABS,0,hb, ps) == OK) {

```


Appendix B-5 - MODIFICATIONS TO MEMORY MANAGER CODE

Jul 6 14:41 1989 ALLOC.C Page 3

```

> 107      /* free old core */
> 108      free_mem(old_base, proc_size);
> 109      /* allocate new core */
> 110      if(hole_base != alloc_mem(proc_size)) {
> 111          sys_lock(RESTORE, &old_state);
> 112          panic("compact - alloc error",1);
> 113      }
> 114
> 115      /* setup process table map */
> 116      rmp->mp_seg[S].mem_phys = rmp->mp_seg[S].mem_phys -
> 117          rmp->mp_seg[T].mem_phys + hole_base;
> 118      rmp->mp_seg[D].mem_phys = rmp->mp_seg[T].mem_len + hole_base;
> 119      rmp->mp_seg[T].mem_phys = hole_base;
> 120      /* tell kernel */
> 121      sys_newmap( (int) (rmp - mproc), rmp->mp_seg, FALSE);
> 122      }
> 123
> 124      }
> 125      sys_lock(RESTORE, &old_state);
> 126      }

```

Appendix B-6 - MODIFICATIONS TO MEMORY MANAGER CODE

Jul 6 14:41 1989 EXEC.C Page 1

```

1  /*=====*/
2  *                                     do_exec                                     *
3  /*=====*/
4  PUBLIC int do_exec()
5  {
6  /* Perform the execve(name, argv, envp) call. The user library builds a
7  * complete stack image, including pointers, args, environ, etc. The stack
8  * is copied to a buffer inside MM, and then to the new core image.
9  */
10
11     char mbuf[MAX_ISTACK_BYTES]; /* buffer for stack and zeroes */
12     char swap_name[4];
13     char *new_sp, *a, *psrc, *pdst;
14     int s, r, in_fd, out_fd, swap_fd, ft, sd;
15     int swapout = FALSE;
16     unsigned loadbytes;
17     vir_bytes src, dst, text_bytes, data_bytes, bss_bytes, stk_bytes, vsp;
18     phys_bytes tot_bytes; /* total space for program, including gap */
19     phys_clicks old_clk, dbuf_len;
20     long sym_bytes, ubytes, xurabytes;
21     vir_clicks sc;
22     struct mproc *mmp, *umppmp;
23     struct stat s_buf, d_buf;
24     struct mem_map tmm[NR_SEGS];
25     union u {
26         char name_buf[MAX_PATH]; /* the name of the file to exec */
27         char zb[ZEROBUF_SIZE]; /* used to zero bss */
28     } u;
29
30     /* Do some validity checks. */
31     mmp = mp;
32     stk_bytes = (vir_bytes) stack_bytes;
33     if (stk_bytes > MAX_ISTACK_BYTES) return(ENOMEM); /* stack too big */
34     if (exec_len <= 0 || exec_len > MAX_PATH) return(EINVAL);
35
36     /* Get the exec file name and see if the file is executable. */
37     src = (vir_bytes) exec_name;
38     dst = (vir_bytes) u.name_buf;
39     r = mem_copy(who, D, (long) src, MM_PROC_NR, D, (long) dst, (long) exec_len);
40     if (r != OK) return(r); /* file name not in user data segment */
41     tell_fs(CHDIR, who, 0, 0); /* temporarily switch to user's directory */
42     in_fd = allowed(u.name_buf, &s_buf, X_BIT); /* is file executable? */
43     tell_fs(CHDIR, 0, 1, 0); /* switch back to MM's own directory */
44     if (in_fd < 0) return(in_fd); /* file was not executable */
45
46     /* Read the file header and extract the segment sizes. */
47     sc = (stk_bytes + CLICK_SIZE - 1) >> CLICK_SHIFT;
48     if (read_header(in_fd, &ft, &text_bytes, &data_bytes, &bss_bytes,
49                     &tot_bytes, &sym_bytes, sc) < 0) {
50         close(in_fd); /* something wrong with header */
51         return(ENOEXEC);
52     }
53

```

Appendix B-7 - MODIFICATIONS TO MEMORY MANAGER CODE

Jul 6 14:41 1989 EXEC.C Page 2

```

54      /* Fetch the stack from the user before destroying the old core image. */
55      src = (vir_bytes) stack_ptr;
56      dst = (vir_bytes) mbuf;
57      if (mem_copy(who, D, (long) src, MM_PROC_NR, D, (long) dst,
58                  (long) stk_bytes) != OK) {
59          close(in_fd);          /* can't fetch stack (e.g. bad virtual addr) */
60          return(EACCES);
61      }
62
63      /* Allocate new memory and release old memory. Fix map and tell kernel. */
64      r = new_mem(text_bytes, data_bytes, bss_bytes, stk_bytes, tot_bytes,
65                  u.zb, ZEROBUF_SIZE, tmm, &old_clk);
66      if (r == EXSWAPD) {
67          swapout = TRUE;
68      } else if (r != OK) {
69          close(in_fd);          /* insufficient core or program too big */
70          return(r);
71      }
72
73      if (!swapout) {
74          /* Patch up stack and copy it from MM to new core image. */
75          vsp = (vir_bytes) (mp->mp_seg[S].mem_vir << CLICK_SHIFT);
76          vsp += (vir_bytes) (mp->mp_seg[S].mem_len << CLICK_SHIFT);
77          vsp -= stk_bytes;
78          patch_ptr(mbuf, vsp);
79          src = (vir_bytes) mbuf;
80          r = mem_copy(MM_PROC_NR, D, (long) src, who, D, (long) vsp,
81                      (long) stk_bytes);
82          if (r != OK) panic("do_exec stack copy err", NO_NUM);
83
84          /* Read in text and data segments. */
85          load_seg(who, in_fd, T, text_bytes);
86          load_seg(who, in_fd, D, data_bytes);
87      #ifdef ATARI_ST
88          if (lseek(in_fd, sym_bytes, 1) < 0)
89              ; /* error */
90          if (relocate(in_fd, mbuf) < 0)
91              ; /* error */
92      #endif
93      } else {
94          /* read T & D from a.out file and write to swap device */
95          /* create swap file */
96          /* change to swap directory */
97          tell_fs(CHDIR, 0, 2, 0);
98          tmpmp = mp;
99          mp = &mproc[MM_PROC_NR];
100
101          /* get swap file name */
102          sd = who;
103          a = swap_name;
104          while(sd) {
105              *a++ = (sd % 10) + 060;
106              sd /= 10;

```

Appendix B-8 - MODIFICATIONS TO MEMORY MANAGER CODE

Jul 6 14:41 1989 EXECC Page 3

```

> 107     }
> 108     *a = 0;
> 109     out_fd = allowed(swap_name, &s_buf, W_BIT); /* swap_file wrtble? */
> 110     s = allowed(".", &d_buf, W_BIT);
> 111     mp = tmpmp;
> 112     if (out_fd >= 0) close(out_fd);
> 113     if (s >= 0) close(s);
> 114     if (s >= 0 && (out_fd >= 0 || out_fd == ENOENT)) {
> 115         /* File is writable or doesn't exist & dir is writable */
> 116         out_fd = creat(swap_name, SWAP_MODE);
> 117     } else {
> 118         tell_fs(CHDIR, 0, 1, 0); /* go back to MM's own dir */
> 119         return(ERROR);
> 120     }
> 121     tell_fs(CHDIR, 0, 1, 0); /* go back to MM's own dir */
> 122     if (out_fd < 0) return(ERROR);
> 123
> 124     /* change existing core to a data buffer */
> 125     /* get max available size */
> 126     dbuf_len = MIN(2047, old_elk);
> 127
> 128     /* setup data buffer */
> 129     rmp->mp_seg[D].mem_phys = rmp->mp_seg[T].mem_phys;
> 130     rmp->mp_seg[D].mem_len = dbuf_len;
> 131
> 132     /* tell kernel about the buffer */
> 133     sys_newmap(who, rmp->mp_seg, TRUE);
> 134
> 135     swap_fd = (who << 8) | (D << 6) | out_fd;
> 136     for(r=0; r<2; r++) {
> 137         tbytes = r ? (long)data_bytes : (long)ext_bytes;
> 138         xtrabytes = (long)
> 139             ((r ? tmm[D].mem_len : tmm[T].mem_len) << CLICK_SHIFT)
> 140             - tbytes;
> 141         while(tbytes) {
> 142             loadbytes = (unsigned)(MIN((long)tbytes,
> 143                                     (long)(dbuf_len << CLICK_SHIFT)));
> 144             /* read from a.out file */
> 145             load_seg(who, in_fd, D, loadbytes);
> 146             /* write it to swap device */
> 147             a = (char *) (rmp->mp_seg[D].mem_vir << CLICK_SHIFT);
> 148
> 149             if(write(swap_fd, a, loadbytes) != loadbytes) {
> 150                 close(in_fd);
> 151                 close(out_fd);
> 152                 panic("do_exec swap device write err", NO_NUM);
> 153             }
> 154             tbytes -= (long)loadbytes;
> 155         }
> 156         if(xtrabytes)
> 157             /* the number of bytes in the swap file for each
> 158             /* segment must be equal to the length of the segment */
> 159             if(lseek(out_fd, xtrabytes, 1) < 0)

```

Appendix B-9 - MODIFICATIONS TO MEMORY MANAGER CODE

Jul 6 14:41 1989 EXEC.C Page 4

```

> 160             printf("lseek error0);
> 161         }
> 162
> 163         /* write stack to swap device */
> 164         /* first fix mproc memory map */
> 165         psrc = (char *)tmm;
> 166         pdst = (char *)rmp->mp_seg;
> 167         while(psrc != ((char *)tmm + sizeof(tmm)))
> 168             *(pdst++) = *(psrc++);
> 169         vsp = (vir_bytes) (rmp->mp_seg[S].mem_vir << CLICK_SHIFT);
> 170         vsp += (vir_bytes) (rmp->mp_seg[S].mem_len << CLICK_SHIFT);
> 171         vsp -= stk_bytes;
> 172         patch_ptr(mbuf, vsp);
> 173         a = mbuf;
> 174         if(lseek(out_fd, (long)(vsp - (tmm[S].mem_vir << CLICK_SHIFT)),1)
> 175             < (long)0)
> 176             printf("lseek error0);
> 177         if(write(out_fd, a, (unsigned)stk_bytes) != stk_bytes) {
> 178             close(in_fd);
> 179             close(out_fd);
> 180             panic("do_exec swap device write err", NO_NUM);
> 181         }
> 182         close(out_fd);
> 183         free_mem(rmp->mp_seg[T].mem_phys, old_clk); /* free the memory */
> 184         rmp->mp_flags |= SWAPPED; /* mark mproc as swapped */
> 185     }
> 186
> 187     close(in_fd); /* don't need exec file any more */
> 188
> 189     /* Take care of setuid/setgid bits. */
> 190     if (s_buf.st_mode & I_SET_UID_BIT) {
> 191         rmp->mp_effuid = s_buf.st_uid;
> 192         tell_fs(SETUID, who, (int) rmp->mp_realuid, (int) rmp->mp_effuid);
> 193     }
> 194     if (s_buf.st_mode & I_SET_GID_BIT) {
> 195         rmp->mp_effgid = s_buf.st_gid;
> 196         tell_fs(SETGID, who, (int) rmp->mp_realgid, (int) rmp->mp_effgid);
> 197     }
> 198
> 199     /* Fix up some 'mproc' fields and tell kernel that exec is done. */
> 200     rmp->mp_deadcild = 0; /* reset swap wait */
> 201     rmp->mp_ssw_map = 0; /* reset all swap signals */
> 202     rmp->mp_catch = 0; /* reset all caught signals */
> 203     rmp->mp_flags &= SEPARATE; /* turn off SEPARATE bit */
> 204     rmp->mp_flags |= ft; /* turn it on for separate I & D files */
> 205
> 206     new_sp = (char *) vsp;
> 207     sys_exec(who, new_sp, swapout);
> 208     return(OK);
> 209 }
> 210
> 211
> 212  */

```

Appendix B-10 - MODIFICATIONS TO MEMORY MANAGER CODE

Jul 6 14:41 1989 EXEC.C Page 5

```

213 *                                     new_mem                                     *
214 *=====*/
215 PRIVATE int new_mem(text_bytes, data_bytes, bss_bytes, stk_bytes,
216                    tot_bytes, bf, zs, trun, old_clk)
217     vir_bytes text_bytes; /* text segment size in bytes */
218     vir_bytes data_bytes; /* size of initialized data in bytes */
219     vir_bytes bss_bytes; /* size of bss in bytes */
220     vir_bytes stk_bytes; /* size of initial stack segment in bytes */
221     phys_bytes tot_bytes; /* total memory to allocate, including gap */
222     char bf[ZEROBUF_SIZE]; /* buffer to use for zeroing data segment */
223     int zs; /* true size of 'bf' */
224     struct mem_map trun[]; /* temporary memory map */
225     phys_clicks *old_clk; /* # of clicks in old process */
226 {
227     /* Allocate new memory and release the old memory. Change the map and report
228      * the new map to the kernel. Zero the new core image's bss, gap and stack.
229      */
230
231     register struct mproc *mmp;
232     vir_clicks text_clicks, data_clicks, gap_clicks, stack_clicks, tot_clicks;
233     phys_clicks new_base;
234     extern phys_clicks alloc_mem();
235     extern phys_clicks tot_hole();
236     #ifdef ATARI_ST
237     phys_clicks base, size;
238     #else
239     char *rzp;
240     vir_bytes vzb;
241     phys_clicks old_clicks;
242     phys_bytes bytes, base, count, bss_offset;
243     #endif
244
245     /* Acquire the new memory. Each of the 4 parts: text, (data+bss), gap,
246      * and stack occupies an integral number of clicks, starting at click
247      * boundary. The data and bss parts are run together with no space.
248      */
249
250     text_clicks = (text_bytes + CLICK_SIZE - 1) >> CLICK_SHIFT;
251     data_clicks = (data_bytes + bss_bytes + CLICK_SIZE - 1) >> CLICK_SHIFT;
252     stack_clicks = (stk_bytes + CLICK_SIZE - 1) >> CLICK_SHIFT;
253     tot_clicks = (tot_bytes + CLICK_SIZE - 1) >> CLICK_SHIFT;
254     gap_clicks = tot_clicks - data_clicks - stack_clicks;
255     if ((int) gap_clicks < 0) return(ENOMEM);
256
257     mmp = mp;
258     #ifdef ATARI_ST
259     old_clicks = (phys_clicks) mmp->mmp_seg[S].mem_len;
260     old_clicks += (mmp->mmp_seg[S].mem_vir - mmp->mmp_seg[D].mem_vir);
261     if (mmp->mmp_flags & SEPARATE) old_clicks += mmp->mmp_seg[T].mem_len;
262     #endif
263
264     /* Check to see if there is a hole big enough. If so, we can risk first
265      * releasing the old core image before allocating the new one, since we

```

Appendix B-11 - MODIFICATIONS TO MEMORY MANAGER CODE

Jul 6 14:41 1989 EXEC.C Page 6

```

266      * know it will succeed. If there is not enough, return failure.
267      */
> 268      if( (text_clicks + tot_clicks) > (tot_hole() + old_clicks) ) {
> 269          /* core is not available, must swap proc out, first get new map */
> 270          tmm[T].mem_len = text_clicks;
> 271          tmm[T].mem_phys = rmp->mp_seg[T].mem_phys;
> 272          tmm[D].mem_len = data_clicks;
> 273          tmm[D].mem_phys = tmm[T].mem_phys + text_clicks;
> 274          tmm[S].mem_len = stack_clicks;
> 275          tmm[S].mem_phys = tmm[D].mem_phys + data_clicks + gap_clicks;
> 276          tmm[T].mem_vir = 0;
> 277          tmm[D].mem_vir = 0;
> 278          tmm[S].mem_vir = tmm[D].mem_vir + data_clicks + gap_clicks;
> 279          *old_clk = old_clicks;
> 280          return(EXSWAPD);
> 281      }
282
283      #ifndef ATARI_ST
284          /* There is enough memory for the new core image. Release the old one. */
285          free_mem(rmp->mp_seg[T].mem_phys, old_clicks); /* free the memory */
286      #endif
287
288      /* We have now passed the point of no return. The old core image has been
289      * forever lost. The call must go through now. Set up and report new map.
290      */
291      new_base = alloc_mem(text_clicks + tot_clicks); /* new core image */
292      if (new_base == NO_MEM) panic("MM hole list is inconsistent", NO_NUM);
293      rmp->mp_seg[T].mem_len = text_clicks;
294      rmp->mp_seg[T].mem_phys = new_base;
295      rmp->mp_seg[D].mem_len = data_clicks;
296      rmp->mp_seg[D].mem_phys = new_base + text_clicks;
297      rmp->mp_seg[S].mem_len = stack_clicks;
298      rmp->mp_seg[S].mem_phys = rmp->mp_seg[D].mem_phys + data_clicks + gap_clicks;
299      #ifdef ATARI_ST
300          rmp->mp_seg[T].mem_vir = rmp->mp_seg[T].mem_phys;
301          rmp->mp_seg[D].mem_vir = rmp->mp_seg[D].mem_phys;
302          rmp->mp_seg[S].mem_vir = rmp->mp_seg[S].mem_phys;
303      #else
304          rmp->mp_seg[T].mem_vir = 0;
305          rmp->mp_seg[D].mem_vir = 0;
306          rmp->mp_seg[S].mem_vir = rmp->mp_seg[D].mem_vir + data_clicks + gap_clicks;
307      #endif
308      #ifdef ATARI_ST
309          sys_fresh(who, rmp->mp_seg, (phys_clicks)(data_bytes >> CLICK_SHIFT),
310                  &base, &size);
311          free_mem(base, size);
312      #else
313          sys_newmap(who, rmp->mp_seg, FALSE); /* report new map to the kernel */
314
315          /* Zero the bss, gap, and stack segment. Start just above text. */
316          for (rzp = &bfb[0]; rzp < &bfb[zs]; rzp++) *rzp = 0; /* clear buffer */
317          bytes = (phys_bytes) (data_clicks + gap_clicks + stack_clicks) << CLICK_SHIFT;
318          vzb = (vir_bytes) bf;

```

Appendix B-12 - MODIFICATIONS TO MEMORY MANAGER CODE

Jul 6 14:41 1989 EXEC.C Page 7

```

319     base = (long) mmp->mp_seg[T].mem_phys + mmp->mp_seg[T].mem_len;
320     base = base << CLICK_SHIFT;
321     base_offset = (data_bytes >> CLICK_SHIFT) << CLICK_SHIFT;
322     base += base_offset;
323     bytes -= base_offset;
324
325     while (bytes > 0) {
326         count = (long) MIN(bytes, (phys_bytes) za);
327         if (mem_copy(MM_PROC_NR, D, (long) vzb, ABS, 0, base, count) != OK)
328             panic("new_mem can't zero", NO_NUM);
329         base += count;
330         bytes -= count;
331     }
332 #endif
333     return(OK);
334 }
335
336
337
338 /*=====
339 *                               load_seg                               *
340 *=====*/
341 PUBLIC load_seg(usr, fd, seg, seg_bytes)
342 int usr; /* user slot in proc table */
343 int fd; /* file descriptor to read from */
344 int seg; /* T or D */
345 vir_bytes seg_bytes; /* how big is the segment */
346 {
347     /* Read in text or data from the exec file and copy to the new core image.
348     * This procedure is a little bit tricky. The logical way to load a segment
349     * would be to read it block by block and copy each block to the user space
350     * one at a time. This is too slow, so we do something dirty here, namely
351     * send the user space and virtual address to the file system in the upper
352     * 10 bits of the file descriptor, and pass it the user virtual address
353     * instead of a MM address. The file system copies the whole segment
354     * directly to user space, bypassing MM completely.
355     */
356
357     int new_fd, bytes;
358     char *ubuf_ptr;
359     struct mproc *mmp;
360
361     new_fd = (usr << 8) | (seg << 6) | fd;
362     mmp = &mproc[usr];
363     ubuf_ptr = (char *) ((vir_bytes)mmp->mp_seg[seg].mem_vir << CLICK_SHIFT);
364     while (seg_bytes) {
365         bytes = 31*1024; /* <= 32767 */
366         if (seg_bytes < bytes)
367             bytes = (int)seg_bytes;
368         if (read(new_fd, ubuf_ptr, bytes) != bytes) {
369             panic("loadseg read err", NO_NUM);
370             break; /* error */
371         }
372         ubuf_ptr += bytes;
373     }

```


Appendix B-13 - MODIFICATIONS TO MEMORY MANAGER CODE

Jul 6 14:41 1989 EXEC.C Page 8

```
372         seg_bytes -= bytes;
373     }
374 }
```

Appendix B-14 - MODIFICATIONS TO MEMORY MANAGER CODE

Jul 6 14:41 1989 FORKEXTT.C Page 1

```

1 /*
2                                     do_fork
3                                     *
4                                     =====
5 PUBLIC int do_fork()
6 {
7 /* The process pointed to by 'mp' has forked. Create a child process. */
8
9 register struct mproc *mmp; /* pointer to parent */
10 register struct mproc *mmc; /* pointer to child */
11
12 int i, child_nr, t;
13 char *sptr, *dptr;
14 phys_clicks prog_clicks, child_base;
15 extern phys_clicks alloc_mem0;
16 extern phys_clicks tot_hole0;
17 int swapout = FALSE;
18 vir_bytes new_sp;
19 #ifndef ATARI_ST
20 long prog_bytes;
21 long parent_abs, child_abs;
22 #endif
23
24 /* If tables might fill up during FORK, don't even start since recovery half
25 * way through is such a nuisance.
26 */
27
28 mmp = mp;
29 if (procs_in_use == NR_PROCS) return(EAGAIN);
30 if (procs_in_use >= NR_PROCS - LAST_FEW && mmp->mp_effuid != 0) return(EAGAIN);
31
32 /* Determine how much memory to allocate. */
33 prog_clicks = (phys_clicks) mmp->mp_seg[S].mem_len;
34 prog_clicks += (mmp->mp_seg[S].mem_vir - mmp->mp_seg[D].mem_vir);
35 #ifndef ATARI_ST
36 if (mmp->mp_flags & SEPARATE) prog_clicks += mmp->mp_seg[T].mem_len;
37 prog_bytes = (long) prog_clicks << CLICK_SHIFT;
38 #endif
39
40 if ( ( (prog_clicks > tot_hole0) ||
41      ( (child_base = alloc_mem(prog_clicks)) == NO_MEM) ) {
42
43     swapout = TRUE;
44     /* adjust parents memory map, if necessary */
45     sys_getsp(who, &new_sp);
46     if (adjust(mmp, (vir_clicks) mmp->mp_seg[D].mem_len, new_sp) != OK)
47         return(EAGAIN);
48 }
49
50 #ifndef ATARI_ST
51 if (! swapout) {
52     /* Create a copy of the parent's core image for the child. */
53     child_abs = (long) child_base << CLICK_SHIFT;
54     parent_abs = (long) mmp->mp_seg[T].mem_phys << CLICK_SHIFT;
55     i = mem_copy(ABS, 0, parent_abs, ABS, 0, child_abs, prog_bytes);
56     if ( i < 0 ) panic("do fork can't copy".i);
57 }
58
59 }
60
61 }
62
63 */

```

Appendix B-15 - MODIFICATIONS TO MEMORY MANAGER CODE

Jul 6 14:41 1989 FORKEXT.C Page 2

```

54     }
55 #endif
56
57 /* Find a slot in 'mproc' for the child process. A slot must exist. */
58 for (rnc = &mproc[0]; rnc < &mproc[NR_PROCS]; rnc++)
59     if ((rnc->mp_flags & IN_USE) == 0) break;
60
61 /* Set up the child and its memory map; copy its 'mproc' slot from parent. */
62 child_nr = (int)(rnc - mproc); /* slot number of the child */
63 spt = (char *) rnc; /* pointer to parent's 'mproc' slot */
64 dptr = (char *) rnc; /* pointer to child's 'mproc' slot */
65 i = sizeof(struct mproc); /* number of bytes in a proc slot. */
66 while (i-- > 0) *dptr++ = *spt++; /* copy from parent slot to child's */
67
68 rnc->mp_parent = who; /* record child's parent */
69 #ifndef ATARI_ST
70 if (!swapout) {
71     rnc->mp_seg[T].mem_phys = child_base;
72     rnc->mp_seg[D].mem_phys = child_base + rnc->mp_seg[T].mem_len;
73     rnc->mp_seg[S].mem_phys = rnc->mp_seg[D].mem_phys +
74         (rnc->mp_seg[S].mem_phys - rnc->mp_seg[D].mem_phys);
75 } else {
76     /* swapout parent's image for child, don't free parent's core */
77     if (swap_out(child_nr, rnc, rnc, FALSE) != OK) {
78         return(EAGAIN);
79     } else {
80         rnc->mp_flags |= FKSWAPPED;
81     }
82 }
83 #endif
84 rnc->mp_exitstatus = 0;
85 rnc->mp_sigstatus = 0;
86 rnc->mp_deadchild = 0; /* reset swap wait */
87 rnc->mp_ssw_map = 0; /* reset all swap signals */
88 procs_in_use++;
89
90 /* Find a free pid for the child and put it in the table. */
91 do {
92     t = 0; /* 't' = 0 means pid still free */
93     next_pid = (next_pid < 30000 ? next_pid + 1 : INIT_PROC_NR + 1);
94     for (mp = &mproc[0]; mp < &mproc[NR_PROCS]; mp++)
95         if (mp->mp_pid == next_pid || mp->mp_progrp == next_pid) {
96             t = 1;
97             break;
98         }
99     rnc->mp_pid = next_pid; /* assign pid to child */
100 } while (t);
101
102 /* Set process group. */
103 if (who == INIT_PROC_NR) rnc->mp_progrp = rnc->mp_pid;
104
105 /* Tell kernel and file system about the (now successful) FORK. */
106 #ifndef ATARI_ST

```

Appendix B-16 - MODIFICATIONS TO MEMORY MANAGER CODE

Jul 6 14:41 1989 FORKEXIT.C Page 3

```

107     sys_fork(who, child_nr, rmc->mp_pid, child_base);
108     #else
> 109     sys_fork(who, child_nr, rmc->mp_pid, swapout);
110     #endif
111
112     tell_fs(FORK, who, child_nr, rmc->mp_pid);
113
114     #ifdef ATARI_ST
115     /* Report child's memory map to kernel. */
116     if(! swapout) {
117         sys_newmap(child_nr, rmc->mp_seg, FALSE);
118     }
119     #endif
120
> 121     if(! swapout) {
122         /* Reply to child to wake it up. */
123         reply(child_nr, 0, 0, NIL_PTR);
124     }
125     return(next_pid); /* child's pid */
126 }
127
128
129 /*=====
130 *                               do_mm_exit                               *
131 *=====*/
132 PUBLIC int do_mm_exit()
133 {
134     /* Perform the exit(status) system call. The real work is done by mm_exit(),
135     * which is also called when a process is killed by a signal.
136     */
137
138     mm_exit(mp, status);
139     dont_reply = TRUE; /* don't reply to newly terminated process */
140     return(OK); /* pro forma return code */
141 }
142
143
144 /*=====
145 *                               mm_exit                               *
146 *=====*/
147 PUBLIC mm_exit(rmp, exit_status)
148 register struct mproc *rmp; /* pointer to the process to be terminated */
149 int exit_status; /* the process' exit status (for parent) */
150 {
151     /* A process is done. If parent is waiting for it, clean it up, else hang. */
152     #ifdef ATARI_ST
153     phys_clicks base, size;
154     #endif
155     phys_clicks s;
156     register int proc_nr = (int)(rmp - inproc);
157
158     /* How to terminate a process is determined by whether or not the
159     * parent process has already done a WAIT. Test to see if it has.

```

Appendix B-17 - MODIFICATIONS TO MEMORY MANAGER CODE

Jul 6 14:41 1989 FORKEXIT.C Page 4

```

160      */
161      rmp->mp_exitstatus = (char) exit_status; /* store status in 'mproc' */
162
163      if (mproc[rmp->mp_parent].mp_flags & WAITING) {
164          if (mproc[rmp->mp_parent].mp_flags & (SWAPPED | FKSWAPPED)) {
165              /* parent is swapped & waiting, delay cleanup by falsely */
166              /* marking child HANGING, setting deadchild, and telling */
167              /* SWAP_TASK */
168              rmp->mp_flags |= HANGING;
169              mproc[rmp->mp_parent].mp_deadchild = proc_nr;
170              mproc[rmp->mp_parent].mp_flags &= WAITING;
171              reply(SWAP_TASK, CORE_IS_NEEDED,
172                  (int)(&mproc[rmp->mp_parent] - mproc), NIL_PTR);
173          } else
174              cleanup(rmp); /* release parent and tell everybody */
175      } else
176          rmp->mp_flags |= HANGING; /* Parent not waiting. Suspend proc */
177
178      /* If the exited process has a timer pending, kill it. */
179      if (rmp->mp_flags & ALARM_ON) set_alarm(proc_nr, (unsigned) 0);
180
181      #ifdef AM_KERNEL
182      /* see if an amoeba transaction was pending or a putrep needed to be done */
183      am_check_sig(proc_nr, 1);
184      #endif
185
186      /* Tell the kernel and FS that the process is no longer runnable. */
187      #ifdef ATARI_ST
188      sys_xit(rmp->mp_parent, proc_nr, &base, &size);
189      free_mem(base, size);
190      #else
191      sys_xit(rmp->mp_parent, proc_nr);
192      #endif
193      tell_fs(EXIT, proc_nr, 0, 0); /* file system can free the proc slot */
194
195      #ifndef ATARI_ST
196      /* Release the memory occupied by the child. */
197      s = (phys_clicks) rmp->mp_seg[S].mem_len;
198      s += (rmp->mp_seg[S].mem_vir - rmp->mp_seg[D].mem_vir);
199      if (rmp->mp_flags & SEPARATE) s += rmp->mp_seg[T].mem_len;
200      free_mem(rmp->mp_seg[T].mem_phys, s); /* free the memory */
201      #endif
202
203      }
204
205      /*=====
206      *
207      * do_wait
208      *=====*/
209      PUBLIC int do_wait()
210      {
211      /* A process wants to wait for a child to terminate. If one is already waiting,
212      * go clean it up and let this WAIT call terminate. Otherwise, really wait.

```

Appendix B-18 - MODIFICATIONS TO MEMORY MANAGER CODE

Jul 6 14:41 1989 FORKEXIT.C Page 5

```

213  */
214
215  register struct mproc *rp;
216  register int children;
217
218  /* A process calling WAIT never gets a reply in the usual way via the
219  * reply() in the main loop. If a child has already exited, the routine
220  * cleanup() sends the reply to awaken the caller.
221  */
222
223  /* Is there a child waiting to be collected? */
224  children = 0;
225  for (rp = &mproc[0]; rp < &mproc[NR_PROCS]; rp++) {
226      if ( (rp->mp_flags & IN_USE) && rp->mp_parent == who) {
227          children++;
228          if (rp->mp_flags & HANGING) {
229              cleanup(rp); /* a child has already exited */
230              dont_reply = TRUE;
231              return(OK);
232          }
233      }
234  }
235
236  /* No child has exited. Wait for one, unless none exists. */
237  if (children > 0) { /* does this process have any children? */
238      mp->mp_flags |= WAITING;
239      dont_reply = TRUE;
240
241      for (rp = &mproc[INIT_PROC_NR + 1]; rp < &mproc[NR_PROCS]; rp++) {
242          if ( (rp->mp_flags & IN_USE) == 0) continue;
243          if ( (rp->mp_flags & (SWAPPED | FKSWAPPED)) &&
244              (rp->mp_flags & (PAUSED | WAITING) == 0) ) {
245              ms.m_source = MM_PROC_NR;
246              ms.m_type = CORE_IS_FREE;
247              send(SWAP_TASK, &ms);
248              break;
249          }
250      }
251      return(OK); /* yes - wait for one to exit */
252  } else
253      return(ECHILD); /* no - parent has no children */
254  }

```

Appendix B-19 - MODIFICATIONS TO MEMORY MANAGER CODE

Jul 6 14:41 1989 MAIN.C Page 1

```

1  /*=====*/
2  *               main               *
3  /*=====*/
4  PUBLIC main()
5  {
6  /* Main routine of the memory manager. */
7
8  int error;
9
10 mm_init();           /* initialize memory manager tables */
11
12 /* This is MM's main loop- get work and do it, forever and forever. */
13 while (TRUE) {
14     /* Wait for message. */
15     get_work();       /* wait for an MM system call */
16     mp = &mproc[who];
17
18     /* Set some flags. */
19     error = OK;
20     dont_reply = FALSE;
21     err_code = -999;
22
23     /* If the call number is valid, perform the call. */
24     if (mm_call < 0 || mm_call >= NCALLS)
25         error = E_BAD_CALL;
26     else
> 27         error = (*call_vec[mm_call])(0);
28
29     /* Send the results back to the user to indicate completion. */
30     if (dont_reply) continue; /* no reply for EXIT and WAIT */
31     if (mm_call == EXEC && error == OK) continue;
32     reply(who, error, result2, res_ptr);
33 }
34 }
35
36
37
38 /*=====*/
39 *               reply               *
40 /*=====*/
41 PUBLIC reply(proc_nr, result, res2, respt)
42 int proc_nr;           /* process to reply to */
43 int result;           /* result of the call (usually OK or error #)*/
44 int res2;             /* secondary result */
45 char *respt;          /* result if pointer */
46 {
47 /* Send a reply to a user process. */
48
49 register struct mproc *proc_ptr;
50
51 /* To make MM robust, check to see if destination is still alive. */
> 52 if(proc_nr != SWAP_TASK) {
53     proc_ptr = &mproc[proc_nr];

```

Appendix B-20 - MODIFICATIONS TO MEMORY MANAGER CODE

Jul 6 14:41 1989 MAIN.C Page 2

```
54     if ( (proc_ptr->mp_flags&IN_USE) == 0 || (proc_ptr->mp_flags&HANGING))
55         return;
56     }
57     reply_type = result;
58     reply_i1 = res2;
59     reply_p1 = respt;
60     if (send(proc_nr, &mm_out) != OK) panic("MM can't reply", NO_NUM);
61 }
```


Appendix B-21 - MODIFICATIONS TO MEMORY MANAGER CODE

Jul 6 14:41 1989 MSWAP.C Page 1

```

> 1  #include "../h/const.h"
> 2  #include "../h/type.h"
> 3  #include "../h/callnr.h"
> 4  #include "../h/com.h"
> 5  #include "../h/error.h"
> 6  #include "../h/stat.h"
> 7  #include "../h/signal.h"
> 8  #include "const.h"
> 9  #include "glo.h"
> 10 #include "mproc.h"
> 11 #include "param.h"
> 12
> 13 #define LONG1          m2_j1  /* message slot to carry long bitmap of */
> 14                               /* pause/wait procs, if proc table has */
> 15                               /* more than 32 slots, another method is */
> 16                               /* needed */
> 17 PRIVATE struct mproc *mp;
> 18
> 19 /*=====
> 20 *                               do_swout                               *
> 21 *=====*/
> 22 do_swout(num)
> 23 int num;
> 24 {
> 25     /* perform request from SWAP_TASK to do swapout of a particular process
> 26     */
> 27     int pnum;
> 28
> 29     mp = &mproc[MM_PROC_NR];    /* mp points to MM */
> 30     if(num)
> 31         pnum = num;
> 32     else
> 33         pnum = mm_in.PROC1;
> 34     if( (pnum < 0) || (pnum > NR_PROCS) ) {
> 35         printf("DSO1: swapout proc out of range: %d0,pnum);
> 36         return(ERROR);
> 37     }
> 38
> 39     rmp = &mproc[pnum];    /* rmp points to swapout proc */
> 40
> 41     if(swap_out(pnum, rmp, rmp, TRUE) == OK) {
> 42         rmp->mp_flags |= SWAPPED;
> 43         dont_reply= TRUE;
> 44         mm_out.m_type = SWAP_OUT_COMPL;
> 45         mm_out.PROC1 = pnum;
> 46         if(rmp->mp_flags & (PAUSED | WAITING))
> 47             mm_out.PROC2 = TRUE;
> 48         else
> 49             mm_out.PROC2 = FALSE;
> 50         if (send(SWAP_TASK,&mm_out) != OK)
> 51             panic("mswap can't send mes", NO_NUM);
> 52         return(SWAP_OUT_COMPL);
> 53     } else {

```

Appendix B-22 - MODIFICATIONS TO MEMORY MANAGER CODE

Jul 6 14:41 1989 MSWAP.C Page 2

```

> 54             result2 = pnum;
> 55             return(SWAP_OUT_FAILED);
> 56         }
> 57     }
> 58
> 59     /*=====
> 60     *                               swap_out                               *
> 61     *=====*/
> 62     swap_out(swap_proc, rmc, rmp, clear_mem)
> 63     int swap_proc;
> 64     struct mproc *rmc; /* child process */
> 65     struct mproc *rmp; /* parent process */
> 66     int clear_mem;
> 67     {
> 68     /* do actual work of swapping out the process image to swap device */
> 69     phys_clicks s;
> 70     int type;
> 71     struct mproc *tmp;
> 72
> 73         tmp = mp;
> 74         mp = &mproc[MM_PROC_NR];
> 75         if(clear_mem) type = 1;
> 76         else type = 2;
> 77
> 78         if(dump_core(rmc, rmp, type) != OK) {
> 79             printf("SWAPOUT ERROR0);
> 80             mp = tmp;
> 81             return(ERROR);
> 82         }
> 83
> 84         mp = tmp;
> 85         if(clear_mem) {
> 86             /* Release the memory occupied by the process. */
> 87             s = (phys_clicks) rmp->mp_seg[S].mem_len;
> 88             s += (rmp->mp_seg[S].mem_vir - rmp->mp_seg[D].mem_vir);
> 89             if (rmp->mp_flags & SEPARATE) s += rmp->mp_seg[T].mem_len;
> 90             free_mem(rmp->mp_seg[T].mem_phys, s); /* free the memory */
> 91         }
> 92         return(OK);
> 93     }
> 94
> 95     /*=====
> 96     *                               do_swin                               *
> 97     *=====*/
> 98     do_swin()
> 99     {
> 100     /* perform request from SWAP_TASK to do swpin of a particular process
> 101     */
> 102     int i, num;
> 103     phys_clicks new_base, tot_clicks;
> 104     extern phys_clicks alloc_mem(), tot_hole();
> 105     extern cleanup();
> 106     struct mproc *pp;

```

Appendix B-23 - MODIFICATIONS TO MEMORY MANAGER CODE

Jul 6 14:41 1989 MSWAP.C Page 3

```

> 107     long bitmap;
> 108
> 109     num = mm_in.PROC1;
> 110     mp = &mproc[MM_PROC_NR]; /* mp points to MM */
> 111     rmp = &mproc[num]; /* rmp points to swapin process */
> 112     tot_clicks = rmp->mp_seg[S].mem_phys -
> 113     rmp->mp_seg[T].mem_phys +
> 114     rmp->mp_seg[S].mem_len;
> 115     if((tot_hole() >= tot_clicks) &&
> 116     ((new_base = alloc_mem(tot_clicks)) != NO_MEM)) {
> 117         rmp->mp_flags &= SWAPPED;
> 118
> 119         if(swap_in(num, rmp, new_base) == OK) {
> 120             if(rmp->mp_flags & FKSWAPPED) {
> 121                 rmp->mp_flags &= FKSWAPPED;
> 122                 /* wake up forked child */
> 123                 reply(num, 0, 0, NIL_PTR);
> 124             }
> 125             /* do all processing required due to */
> 126             /* proc being swapped out, messy stuff */
> 127             /* process signals recvd during swapout */
> 128             if(rmp->mp_ssw_map) {
> 129                 for(i=1; i<=NR_SIGS; i++) {
> 130                     if(rmp->mp_ssw_map & (1 << i-1))
> 131                         sig_proc(rmp, i);
> 132                 }
> 133                 rmp->mp_ssw_map = 0;
> 134             }
> 135             /* if proc was WAIT & SWAPPED & a child died */
> 136             if(rmp->mp_deadchild) {
> 137                 cleanup(&mproc[rmp->mp_deadchild]);
> 138                 rmp->mp_deadchild = 0;
> 139             }
> 140             /* wake up P/W proc awakened by signal while SWAPPED */
> 141             if(rmp->mp_flags & WASPWS) {
> 142                 reply(num, EINTR, 0, NIL_PTR);
> 143                 rmp->mp_flags &= WASPWS;
> 144             }
> 145             result2 = num;
> 146             return(SWAP_IN_COMPL);
> 147         } else {
> 148             panic("SWAPIN ERROR", NO_NUM);
> 149         }
> 150     } else {
> 151         /* no core available, send kernel list of PAUSED & */
> 152         /* WAITING procs and amount of core needed */
> 153         rmp->mp_flags |= SWAPPED;
> 154         bitmap = 0;
> 155         for(rpp = &mproc[INIT_PROC_NR + 1];
> 156             rpp < &mproc[NR_PROCS]; rpp++) {
> 157             if(((rpp->mp_flags & IN_USE) == 0) ||
> 158                 (rpp->mp_flags & (HANGING | SWAPPED | FKSWAPPED)))
> 159                 continue;

```

Appendix B-24 - MODIFICATIONS TO MEMORY MANAGER CODE

Jul 6 14:41 1989 MSWAP.C Page 4

```

> 160             if(rpp->mp_flags & (PAUSED | WAITING))
> 161                 bitmap |= 1 << (int)(rpp - mproc);
> 162             }
> 163             mm_out_m_type = SWAP_IN_FAILED;
> 164             mm_out.PROC1 = (unsigned)(tot_clicks - tot_hole0);
> 165             mm_out.LONG1 = bitmap;
> 166             if (sendrec(SWAP_TASK,&mm_out) != OK)
> 167                 panic("mswap can't send mes", NO_NUM);
> 168             if(mm_out_m_type != SWAP_OUT_REQ) {
> 169                 dont_reply = TRUE;
> 170                 return(SWAP_IN_FAILED);
> 171             } else {
> 172                 return(do_swait(mm_out.PROC1));
> 173             }
> 174         }
> 175     }
> 176
> 177
> 178     /* ===== swap_in ===== */
> 179     *
> 180     /* ===== */
> 181     PUBLIC int swap_in(swap_proc, mmp, new_base)
> 182     int swap_proc;
> 183     struct mproc *mmp;
> 184     phys_clicks new_base;
> 185     {
> 186     /* do actual work of swapping in the process image from swap device */
> 187
> 188     int fd, tmp;
> 189     char name_buf[5]; /* the name of the file to swapin */
> 190     char zb[ZEROBUF_SIZE]; /* used to zero bss */
> 191     struct stat s_buf;
> 192     phys_clicks gap_clicks;
> 193     extern int load_seg0;
> 194     char *rzp;
> 195     vir_bytes vzb;
> 196     phys_bytes bytes, base, count;
> 197
> 198     /* Do some validity checks. */
> 199
> 200     /* Get the swap file name */
> 201     rzp = name_buf;
> 202     tmp = swap_proc;
> 203     while(tmp) { /* same algorithm used in dump_core to determine swapname */
> 204         *rzp++ = (tmp % 10) + 060;
> 205         tmp /= 10;
> 206     }
> 207     *rzp = 0;
> 208
> 209     tell_fs(CHDIR, 0, 2, 0); /* temporarily switch to swap dir */
> 210     fd = allowed(name_buf, &s_buf, R_BIT); /* is file readable? */
> 211     tell_fs(CHDIR, 0, 1, 0); /* switch back to MM's own directory */
> 212     if (fd < 0) {

```

Appendix B-25 - MODIFICATIONS TO MEMORY MANAGER CODE

Jul 6 14:41 1989 MSWAP.C Page 5

```

> 213         printf("SWAP FILE not readable");
> 214         return(ERROR);      /* file was not readable */
> 215     }
> 216
> 217     /* Fix map with new memory and tell kernel. */
> 218     gap_clicks = rmp->mp_seg[S].mem_phys -
> 219         (rmp->mp_seg[D].mem_phys +
> 220          rmp->mp_seg[D].mem_len);
> 221     rmp->mp_seg[T].mem_phys = new_base;
> 222     rmp->mp_seg[D].mem_phys = new_base + rmp->mp_seg[T].mem_len;
> 223     rmp->mp_seg[S].mem_phys = rmp->mp_seg[D].mem_phys +
> 224         rmp->mp_seg[D].mem_len +
> 225         gap_clicks;
> 226     rmp->mp_seg[T].mem_vir = 0;
> 227     rmp->mp_seg[D].mem_vir = 0;
> 228     rmp->mp_seg[S].mem_vir = rmp->mp_seg[D].mem_vir +
> 229         rmp->mp_seg[D].mem_len +
> 230         gap_clicks;
> 231     sys_newmap(swap_proc, rmp->mp_seg, FALSE); /* report new map to the kernel */
> 232
> 233     /* Zero the gap */
> 234     for (rzb = zb; rzb < &z[ZEROBUF_SIZE]; rzb++) *rzb = 0; /* clear buffer */
> 235     bytes = (phys_bytes) gap_clicks << CLICK_SHIFT;
> 236     vzb = (vir_bytes) zb;
> 237     base = (long) rmp->mp_seg[D].mem_phys + rmp->mp_seg[D].mem_len;
> 238     base = base << CLICK_SHIFT;
> 239
> 240     while (bytes > 0) {
> 241         count = (long) MIN(bytes, (phys_bytes) ZEROBUF_SIZE);
> 242         if (mem_copy(MM_PROC_NR, D, (long) vzb, ABS, 0, base, count) != OK)
> 243             panic("new_mem can't zero", NO_NUM);
> 244         base += count;
> 245         bytes -= count;
> 246     }
> 247
> 248     /* Read in text, data and stack segments. */
> 249     load_seg(swap_proc, fd, T, (phys_bytes) rmp->mp_seg[T].mem_len << CLICK_SHIFT);
> 250     load_seg(swap_proc, fd, D, (phys_bytes) rmp->mp_seg[D].mem_len << CLICK_SHIFT);
> 251     load_seg(swap_proc, fd, S, (phys_bytes) rmp->mp_seg[S].mem_len << CLICK_SHIFT);
> 252     close(fd); /* don't need swap file any more */
> 253
> 254     tell_fs(CHDIR, 0, 2, 0); /* temporarily switch to swap dir */
> 255     if (unlink(name_buf) != 0)
> 256         printf("unlink error: %s0,name_buf);
> 257     tell_fs(CHDIR, 0, 1, 0); /* switch back to MM's own directory */
> 258
> 259     return(OK);
> 260 }

```

Appendix B-26 - MODIFICATIONS TO MEMORY MANAGER CODE

Jul 6 14:41 1989 SIGNALC Page 1

```

1  /*=====
2  *                      check_sig                      *
3  *=====*/
4  PRIVATE int check_sig(proc_id, sig_nr, send_uid)
5  int proc_id;          /* pid of process to signal, or 0 or -1 */
6  int sig_nr;           /* which signal to send (1-16) */
7  uid send_uid;         /* identity of process sending the signal */
8  {
9      /* Check to see if it is possible to send a signal. The signal may have to be
10     * sent to a group of processes. This routine is invoked by the KILL system
11     * call, and also when the kernel catches a DEL or other signal. SIGALRM too.
12     */
13
14     register struct mproc *mproc;
15     int count, send_sig;
16     unshort mask;
17     extern unshort core_bits;
18
19     if (sig_nr < 1 || sig_nr > NR_SIGS) return(EINVAL);
20     count = 0;          /* count # of signals sent */
21     mask = 1 << (sig_nr - 1);
22
23     /* Search the proc table for processes to signal. Several tests are made:
24     *   - if proc's uid != sender's, and sender is not superuser, don't signal
25     *   - if specific process requested (i.e., "procpid" > 0), check for match
26     *   - if a process has already exited, it can't receive signals
27     *   - if "proc_id" is 0 signal everyone in same process group except caller
28     */
29     for (mproc = &mproc[INIT_PROC_NR + 1]; mproc < &mproc[NR_PROCS]; mproc++) {
30         if ((mproc->mp_flags & IN_USE) == 0) continue;
31         send_sig = TRUE; /* if it's FALSE at end of loop, don't signal */
32         if (send_uid != mproc->mp_effuid && send_uid != SUPER_USER) send_sig = FALSE;
33         if (proc_id > 0 && proc_id != mproc->mp_pid) send_sig = FALSE;
34         if (mproc->mp_flags & HANGING) send_sig = FALSE; /*don't wake the dead*/
35         if (proc_id == 0 && mproc->mp_procgrp != mproc->mp_procgrp) send_sig = FALSE;
36         if (send_uid == SUPER_USER && proc_id == -1) send_sig = TRUE;
37
38         /* SIGALARM is a little special. When a process exits, a clock signal
39         * can arrive just as the timer is being turned off. Also, turn off
40         * ALARM_ON bit when timer goes off to keep it accurate.
41         */
42         if (sig_nr == SIGALRM) {
43             if ((mproc->mp_flags & ALARM_ON) == 0) continue;
44             if (send_sig) mproc->mp_flags &= ALARM_ON;
45         }
46
47         if (send_sig == FALSE) continue;
48         count++;
49         if (mproc->mp_ignore & mask) continue;
50
51     #ifdef AM_KERNEL
52         /* see if an amoeba transaction should be signalled */
53         Tfs = am_check_sig(mproc - mproc, 0);

```

Appendix B-27 - MODIFICATIONS TO MEMORY MANAGER CODE

Jul 6 14:41 1989 SIGNALC Page 2

```

54 #endif
55
> 56 /* Send the signal or kill the process, possibly with core dump. */
> 57 if(mmp->mp_flags & (SWAPPED | FKSWAPPED))
> 58     /* proc is SWAPPED, delay processing until swapped in */
> 59     mmp->mp_ssw_map |= 1 << (sig_nr - 1);
> 60 else
> 61     sig_proc(mmp, sig_nr);
62
63 /* If process is hanging on PAUSE, WAIT, tty, pipe, etc. release it. */
64 unpause((int)(mmp - mproc)); /* check to see if process is paused */
65 if (proc_id > 0) break; /* only one process being signaled */
66 }
67
68 /* If the calling process has killed itself, don't reply. */
69 if ((mmp->mp_flags & IN_USE) == 0 || (mmp->mp_flags & HANGING)) dont_reply = TRUE;
70 return(count > 0 ? OK : ESRCH);
71 }
72
73
74 /*=====
75 * do_pause
76 *=====*/
77 PUBLIC int do_pause()
78 {
79     /* Perform the pause() system call. */
80     struct mproc *mmp;
81
82     mmp->mp_flags |= PAUSED; /* turn on PAUSE bit */
83     dont_reply = TRUE;
84
85     for (mmp = &mproc[INIT_PROC_NR + 1]; mmp < &mproc[NR_PROCS]; mmp++) {
86         if ((mmp->mp_flags & IN_USE) == 0) continue;
87         if ((mmp->mp_flags & (SWAPPED | FKSWAPPED)) &&
88             (mmp->mp_flags & (PAUSED | WAITING) == 0)) {
89             reply(SWAP_TASK, CORE_IS_FREE, 0, NIL_PTR);
90             break;
91         }
92     }
93     return(OK);
94 }
95
96
97 /*=====
98 * unpause
99 *=====*/
100 PRIVATE unpause(pro)
101 int pro; /* which process number */
102 {
103     /* A signal is to be sent to a process. If that process is hanging on a
104     * system call, the system call must be terminated with EINTR. Possible
105     * calls are PAUSE, WAIT, READ and WRITE, the latter two for pipes and ttys.
106     * First check if the process is hanging on PAUSE or WAIT. If not, tell FS,

```

Appendix B-28 - MODIFICATIONS TO MEMORY MANAGER CODE

Jul 6 14:41 1989 SIGNALC Page 3

```

107  * so it can check for READs and WRITEs from pipes, tty's and the like.
108  */
109
110  register struct mproc *rmp;
111
112  rmp = &mproc[pro];
113
114  /* Check to see if process is hanging on a PAUSE call. */
115  if ((rmp->mp_flags & PAUSED) && (rmp->mp_flags & HANGING) == 0) {
116      rmp->mp_flags &= PAUSED; /* turn off PAUSED bit */
117      if((rmp->mp_flags & (SWAPPED | FKSWAPPED))) {
118          /* PAUSED swapped proc awakened, delay notice & tell kernel */
119          rmp->mp_flags |= WASPWS;
120          reply(SWAP_TASK, CORE_IS_NEEDED, pro, NIL_PTR);
121      } else
122          reply(pro, EINTR, 0, NIL_PTR);
123      return;
124  }
125
126  /* Check to see if process is hanging on a WAIT call. */
127  if ((rmp->mp_flags & WAITING) && (rmp->mp_flags & HANGING) == 0) {
128      rmp->mp_flags &= WAITING; /* turn off WAITING bit */
129      if((rmp->mp_flags & (SWAPPED | FKSWAPPED))) {
130          /* PAUSED swapped proc awakened, delay notice & tell kernel */
131          rmp->mp_flags |= WASPWS;
132          reply(SWAP_TASK, CORE_IS_NEEDED, pro, NIL_PTR);
133      } else
134          reply(pro, EINTR, 0, NIL_PTR);
135      return;
136  }
137
138  #ifdef AM_KERNEL
139      /* if it was an amoeba transaction, it is already tidied up by now. */
140      if (TIFs)
141          #endif
142      /* Process is not hanging on an MM call. Ask FS to take a look. */
143      tell_fs(UNPAUSE, pro, 0, 0);
144
145      return;
146  }
147
148
149  /*=====
150  *                               *
151  *                               *
152  *=====*/
153  PUBLIC dump_core(rmc, rmp, type)
154  struct mproc *rmc; /* child proc for swapout */
155  struct mproc *rmp; /* whose core is to be dumped */
156  int type; /* 0-dump core; 1-swapout; 2-sp_out w/oadjust */
157  {
158      /* Make a core dump on the file "core", if possible. */
159      struct stat s_buf, d_buf;

```


Appendix B-29 - MODIFICATIONS TO MEMORY MANAGER CODE

Jul 6 14:41 1989 SIGNAL.C Page 4

```

> 160     int i, r, s, new_fd, slot, dir, flag, bytes;
> 161     vir_bytes v_buf, c, new_sp;
> 162     char *a;
> 163     struct mproc *xmp;
> 164     extern char core_name[];
> 165     extern adjust();
> 166     char swap_name[4];
> 167
> 168     slot = (int)(rme - mproc);
> 169     /* Change to working directory of dumper. */
> 170     if (type > 0) {
> 171         dir = 0; /* swapout */
> 172         flag = 2;
> 173     } else {
> 174         dir = slot; /* dump_core */
> 175         flag = 0;
> 176     }
> 177     tell_fs(CHDIR, dir, flag, 0);
> 178
> 179     /* Can file be written? */
> 180     if ( (type == 0) && (rme->mp_realuid != rme->mp_effuid) ) {
> 181         tell_fs(CHDIR, 0, 1, 0); /* go back to MM's directory */
> 182         return(ERROR);
> 183     }
> 184     if (type == 0) {
> 185         xmp = mp; /* allowed() looks at 'mp' */
> 186         rme = rmc;
> 187     }
> 188
> 189     if (type > 0) {
> 190         dir = slot;
> 191         a = swap_name;
> 192         while (dir) {
> 193             *a++ = (dir % 10) + 0x60;
> 194             dir /= 10;
> 195         }
> 196         *a = 0;
> 197         r = allowed(swap_name, &s_buf, W_BIT); /* is swap_file writable */
> 198     } else
> 199         r = allowed(core_name, &s_buf, W_BIT); /* is core_file writable */
> 200
> 201     s = allowed(".", &d_buf, W_BIT); /* is directory writable? */
> 202
> 203     if (type == 0)
> 204         mp = xmp;
> 205     if (r >= 0) close(r);
> 206     if (s >= 0) close(s);
> 207     if ((type > 0) || (rme->mp_effuid == SUPER_USER))
> 208         r = 0; /* su can always dump core */
> 209
> 210     if (s >= 0 && (r >= 0 || r == ENOENT)) {
> 211         /* Either file is writable or it doesn't exist & dir is writable */
> 212         if (type > 0)

```

Appendix B-30 - MODIFICATIONS TO MEMORY MANAGER CODE

Jul 6 14:41 1989 SIGNALC Page 5

```

> 213         r = creat(swap_name, SWAP_MODE);
> 214     else
> 215         r = creat(core_name, CORE_MODE);
> 216     tell_fst(CHDIR, 0, 1, 0); /* go back to MM's own dir */
> 217     if (r < 0) {
> 218         printf("create error0);
> 219         return(ERROR);
> 220     }
> 221
> 222     if(type != 2) {
> 223         /* adjust memory map, if necessary (already done for type 2) */
> 224         sys_getsp(slot, &new_sp);
> 225         if(adjust(mmc, (vir_clicks) mmc->mp_seg[D].mem_len, new_sp)
> 226             != OK) {
> 227             printf("ADJUST ERROR0);
> 228             close(r);
> 229             return(ERROR);
> 230         }
> 231     }
> 232
> 233     if(type == 0) {
> 234         mmc->mp_sigstatus |= DUMPED;
> 235
> 236         /* First write the memory map of all segments on core file. */
> 237         if (write(r, (char *) mmc->mp_seg, sizeof(mmc->mp_seg)) < 0) {
> 238             close(r);
> 239             printf("write memory map error0);
> 240             return(ERROR);
> 241         }
> 242     }
> 243
> 244     if(type == 2)
> 245         slot = (int) (mmp - mproc);
> 246     /* Now loop through segments and write the segments themselves out. */
> 247     for (i = 0; i < NR_SEGS; i++) {
> 248         a = (char *) (mmc->mp_seg[i].mem_vir << CLICK_SHIFT);
> 249         c = (int) (mmc->mp_seg[i].mem_len << CLICK_SHIFT);
> 250         new_fd = ((int)(mmp - mproc) << 8) | (i << 6) | r;
> 251
> 252         /* Dump segment. */
> 253         while(c) {
> 254             bytes = 31 * 1024;
> 255             if(c < (vir_bytes)bytes)
> 256                 bytes = (int)c;
> 257             if (write(new_fd, a, bytes) != bytes) {
> 258                 close(r);
> 259                 printf("write segment error0);
> 260                 return(ERROR);
> 261             }
> 262             a += bytes;
> 263             c -= bytes;
> 264         }
> 265     }

```

Appendix B-31 - MODIFICATIONS TO MEMORY MANAGER CODE

Jul 6 14:41 1989 SIGNALC Page 6

```
> 266     } else {
> 267         printf("swap file or dir is not writable");
> 268         tell_fs(CHDIR, 0, 1, 0);      /* go back to MM's own dir */
> 269         close(r);
> 270         return(ERROR);
> 271     }
> 272
> 273     close(r);
> 274     return(OK);
> 275 }
```

Appendix B-32 - MODIFICATIONS TO MEMORY MANAGER CODE

Jul 6 14:41 1989 TABLE.C Page 1

```

1  /* This file contains the table used to map system call numbers onto the
2  * routines that perform them.
3  */
4
5  #include "../h/const.h"
6  #include "../h/type.h"
7  #include "const.h"
8
9  #undef EXTERN
10 #define EXTERN
11
12 #include "../h/callnr.h"
13 #include "glo.h"
14 #include "mproc.h"
15 #include "param.h"
16
17 /* Miscellaneous */
18 char core_name[] = {"core"}; /* file name where core images are produced */
> 19 char swap_name[] = {"xxx"}; /* swap device file names */
20 #ifdef ATARI_ST
21 /*
22  * Creating core files is disabled, except for SIGQUIT and SIGIOT.
23  * Set core_bits to 0x0EFC if you want compatibility with UNIX V7.
24  */
25 ushort core_bits = 0x0EFC; /* which signals cause core images */
26 #else
27 ushort core_bits = 0x0EFC; /* which signals cause core images */
28 #endif
29
30 extern char mm_stack[];
31 char *stackpt = &mm_stack[MM_STACK_BYTES]; /* initial stack pointer */
32
33 extern do_mm_exit(), do_fork(), do_wait(), do_brk(), do_getset(), do_exec();
34 extern do_signal(), do_kill(), do_pause(), do_alarm();
> 35 extern no_sys(), do_ksig(), do_brk2(), do_swin(), do_swout();
36
37 #ifdef AM_KERNEL
38 extern do_amoeba();
39 #endif
40
41 int (*call_vec[NCALLS])() = {
42     no_sys, /* 0 = unused */
43     do_mm_exit, /* 1 = exit */
44     do_fork, /* 2 = fork */
45     no_sys, /* 3 = read */
46     no_sys, /* 4 = write */
47     no_sys, /* 5 = open */
48     no_sys, /* 6 = close */
49     do_wait, /* 7 = wait */
50     no_sys, /* 8 = creat */
51     no_sys, /* 9 = link */
52     no_sys, /* 10 = unlink */
53     no_sys, /* 11 = exec */

```

Appendix B-33 - MODIFICATIONS TO MEMORY MANAGER CODE

Jul 6 14:41 1989 TABLE.C Page 2

```

54      no_sys,      /* 12 = chdir      */
55      no_sys,      /* 13 = time      */
56      no_sys,      /* 14 = mknod     */
57      no_sys,      /* 15 = chmod     */
58      no_sys,      /* 16 = chown     */
59      do_brk,      /* 17 = break     */
60      no_sys,      /* 18 = stat      */
61      no_sys,      /* 19 = lseek     */
62      do_getset, /* 20 = getpid    */
63      no_sys,      /* 21 = mount     */
64      no_sys,      /* 22 = umount    */
65      do_getset, /* 23 = setuid    */
66      do_getset, /* 24 = setuid    */
67      no_sys,      /* 25 = stime     */
68      no_sys,      /* 26 = (ptrace) */
69      do_alarm, /* 27 = alarm     */
70      no_sys,      /* 28 = fstat     */
71      do_pause, /* 29 = pause     */
72      no_sys,      /* 30 = utime     */
73      no_sys,      /* 31 = (stry)    */
74      no_sys,      /* 32 = (gty)     */
75      no_sys,      /* 33 = access    */
76      no_sys,      /* 34 = (nice)    */
77      no_sys,      /* 35 = (ftime)   */
78      no_sys,      /* 36 = sync      */
79      do_kill, /* 37 = kill      */
80      no_sys,      /* 38 = unused    */
81      no_sys,      /* 39 = unused    */
82      no_sys,      /* 40 = unused    */
83      no_sys,      /* 41 = dup       */
84      no_sys,      /* 42 = pipe      */
85      no_sys,      /* 43 = times     */
86      no_sys,      /* 44 = (prof)    */
87      no_sys,      /* 45 = unused    */
88      do_getset, /* 46 = setgid    */
89      do_getset, /* 47 = getgid    */
90      do_signal, /* 48 = sig*      */
91      no_sys,      /* 49 = unused    */
92      no_sys,      /* 50 = unused    */
93      no_sys,      /* 51 = (acct)    */
94      no_sys,      /* 52 = (phys)    */
95      no_sys,      /* 53 = (lock)    */
96      no_sys,      /* 54 = ioctl     */
97      no_sys,      /* 55 = unused    */
98      no_sys,      /* 56 = (mpx)     */
99      no_sys,      /* 57 = unused    */
100     no_sys,      /* 58 = unused    */
101     do_exec, /* 59 = exec*     */
102     no_sys,      /* 60 = umask     */
103     no_sys,      /* 61 = chroot    */
104     no_sys,      /* 62 = unused    */
105     no_sys,      /* 63 = unused    */
106

```

Appendix B-34 - MODIFICATIONS TO MEMORY MANAGER CODE

Jul 6 14:41 1989 TABLE.C Page 3

```

107         do_ksig, /* 64 = KSIG: signals originating in the kernel */
108         no_sys, /* 65 = UNPAUSE */
109         do_brk2, /* 66 = BRK2 (used to tell MM size of FS_INIT) */
110         no_sys, /* 67 = REVIVE */
111         no_sys, /* 68 = TASK_REPLY */
112     #ifdef i8088
113     #ifdef AM_KERNEL
114         do_amoeba, /* 69 = AMOEBA SYSTEM CALL */
115     #else
116         no_sys, /* 69 = AMOEBA SYSTEM CALL */
117     #endif
118     #endif i8088
119     #ifdef SWAPER
> 120         do_swin, /* 70 = swap in */
> 121         do_swout, /* 71 = swap out */
122     #else
123         no_sys, /* 70 = swap in */
124         no_sys, /* 71 = swap out */
125     #endif
126 };

```

APPENDIX C - MODIFICATIONS TO FILE SYSTEM CODE

Appendix C-2 - MODIFICATIONS TO FILE SYSTEM CODE

Jul 6 11:32 1989 MAIN.C Page 1

```

1  /*=====
2  *                               fs_init                               *
3  *=====*/
4  PRIVATE fs_init()
5  {
6  /* Initialize global variables, tables, etc. */
7
8      register struct inode *rip;
9      int i;
10     extern struct inode *get_inode(), *swap_node;
11
12     buf_pool();                /* initialize buffer pool */
13     load_ram();                /* Load RAM disk from root diskette. */
14     load_super();              /* Load super block for root device */
15
16     /* Initialize the 'fproc' fields for process 0 and process 2. */
17     for (i = 0; i < 3; i+= 2) {
18         fp = &fproc[i];
19         rip = get_inode(ROOT_DEV, ROOT_INODE);
20         fp->fp_rootdir = rip;
21         dup_inode(rip);
22         fp->fp_workdir = rip;
23         fp->fp_realluid = (uid) SYS_UID;
24         fp->fp_effuid = (uid) SYS_UID;
25         fp->fp_realgid = (gid) SYS_GID;
26         fp->fp_effgid = (gid) SYS_GID;
27         fp->fp_umask = 0;
28     }
29     swap_node = NIL_INODE;
30
31     /* Certain relations must hold for the file system to work at all. */
32     if (ZONE_NUM_SIZE != 2) panic("ZONE_NUM_SIZE != 2", NO_NUM);
33     if (SUPER_SIZE > BLOCK_SIZE) panic("SUPER_SIZE > BLOCK_SIZE", NO_NUM);
34     if (BLOCK_SIZE % INODE_SIZE != 0) panic("BLOCK_SIZE % INODE_SIZE != 0", NO_NUM);
35     if (NR_FDS > 127) panic("NR_FDS > 127", NO_NUM);
36     if (NR_BUFS < 6) panic("NR_BUFS < 6", NO_NUM);
37     if (sizeof(d_inode) != 32) panic("inode size != 32", NO_NUM);
38 }

```


Appendix C-3 - MODIFICATIONS TO FILE SYSTEM CODE

Jul 6 11:33 1989 STADIRC Page 1

```

1  /* This file contains the code for performing four system calls relating to
2  * status and directories.
3  *
4  * The entry points into this file are
5  * do_chdir:    perform the CHDIR system call
6  * do_chroot:   perform the CHROOT system call
7  * do_stat:     perform the STAT system call
8  * do_fstat:    perform the FSTAT system call
9  */
10
11 #include "../n/const.h"
12 #include "../n/type.h"
13 #include "../n/error.h"
14 #include "../n/stat.h"
15 #include "const.h"
16 #include "type.h"
17 #include "file.h"
18 #include "fproc.h"
19 #include "glo.h"
20 #include "inode.h"
21 #include "param.h"
22 extern struct inode *swap_node;
23 char swap_dir[] = {"usr/swap"};
24
25 /* =====
26 *                               do_chdir                               *
27 * ===== */
28 PUBLIC int do_chdir()
29 {
30     /* Change directory. This function is also called by MM to simulate a chdir
31     * in order to do EXEC, etc.
32     */
33
34     register struct fproc *rfp;
35     int r;
36
37     if (who == MM_PROC_NR) {
38         if (cd_flag == 2) {
39             if (swap_node == NIL_INODE)
40                 if ((r=change(&swap_node, swap_dir, 9)) != OK) {
41                     return(r);
42                 }
43             dup_inode(swap_node);
44             put_inode(fp->fp_workdir);
45             fp->fp_workdir = swap_node;
46             fp->fp_effuid = SUPER_USER;
47             return(OK);
48         }
49     } else {
50         rfp = &fproc[slot1];
51         put_inode(fp->fp_workdir);
52         fp->fp_workdir = (cd_flag == 1 ? fp->fp_rootdir : rfp->fp_workdir);
53         dup_inode(fp->fp_workdir);

```

Appendix C-4 - MODIFICATIONS TO FILE SYSTEM CODE

Jul 6 11:33 1989 STADIR.C Page 2

```

54         fp->fp_effuid = (cd_flag == 1 ? SUPER_USER : rfp->fp_effuid);
55         return(OK);
56     }
57 }
58
59 /* Perform the chdir(name) system call. */
60 return change(&fp->fp_workdir, name, name_length);
61 }
62
63 /* =====
64  *                      change                      *
65  * ===== */
66 PRIVATE int change(iip, name_ptr, len)
67 struct inode **iip;          /* pointer to the inode pointer for the dir */
68 char *name_ptr;              /* pointer to the directory name to change to */
69 int len;                      /* length of the directory name string */
70                               /* if == 0, then use name_ptr directly */
71 {
72     /* Do the actual work for chdir() and chroot(). */
73
74     struct inode *rip;
75     register int r;
76     extern struct inode *cat_path();
77
78     /* Try to open the new directory. */
79     r = 0;
80     if(cd_flag == 2)
81     {
82         user_path[r] = name_ptr[r];
83     } while (name_ptr[r++] != 0);
84
85     else
86     {
87         if (fetch_name(name_ptr, len, M3) != OK) return(err_code);
88
89         if ((rip = cat_path(user_path)) == NIL_INODE) return(err_code);
90         /* It must be a directory and also be searchable. */
91         if ((rip->i_mode & I_TYPE) != I_DIRECTORY)
92             r = ENOTDIR;
93         else
94             r = forbidden(rip, X_BIT, 0); /* check if dir is searchable */
95
96         /* If error, return inode. */
97         if (r != OK) {
98             put_inode(rip);
99             return(r);
100         }
101
102         /* Everything is OK. Make the change. */
103         put_inode(*iip); /* release the old directory */
104         *iip = rip;      /* acquire the new one */
105         return(OK);
106     }

```

An Implementation of Process Swapping in MINIX
(A Message Passing Oriented Operating System)

by

Stanley George Kobylanski

B.S., Pennsylvania State University, 1972

AN ABSTRACT OF A MASTER'S REPORT

submitted in partial fulfillment of the

requirements for the degree

MASTER OF SCIENCE

Department of Computing and Information Sciences

KANSAS STATE UNIVERSITY
Manhattan, Kansas

1989

ABSTRACT

MINIX is a small, general purpose, client-server (transaction oriented) operating system that runs on IBM-PC compatible computers. It provides a UNIX, version 7 based interface and includes many of the standard UNIX support and utility programs. It's intended use is to provide a vehicle for teaching operating system concepts. A limitation of the system is the small amount of main memory available for user processes. This paper provides a solution to that problem by describing an implementation of process swapping.